



London Public Library
REFERENCE



Digitized by the Internet Archive
in 2012

<http://archive.org/details/newencyclopaedia23ency>

The New Encyclopædia Britannica

Volume 23

MACROPÆDIA

Knowledge in Depth

FOUNDED 1768

15TH EDITION



Encyclopædia Britannica, Inc.
Jacob E. Safra, Chairman of the Board
Ilan Yeshua, Chief Executive Officer

Chicago
London/New Delhi/Paris/Seoul
Sydney/Taipei/Tokyo

First Edition 1768-1771
Second Edition 1777-1784
Third Edition 1788-1797
Supplement 1801
Fourth Edition 1801-1809
Fifth Edition 1815
Sixth Edition 1820-1823
Supplement 1815-1824
Seventh Edition 1830-1842
Eighth Edition 1852-1860
Ninth Edition 1875-1889
Tenth Edition 1902-1903

Eleventh Edition

© 1911

By Encyclopædia Britannica, Inc.

Twelfth Edition

© 1922

By Encyclopædia Britannica, Inc.

Thirteenth Edition

© 1926

By Encyclopædia Britannica, Inc.

Fourteenth Edition

© 1929, 1930, 1932, 1933, 1936, 1937, 1938, 1939, 1940, 1941, 1942, 1943,
1944, 1945, 1946, 1947, 1948, 1949, 1950, 1951, 1952, 1953, 1954,
1955, 1956, 1957, 1958, 1959, 1960, 1961, 1962, 1963, 1964,
1965, 1966, 1967, 1968, 1969, 1970, 1971, 1972, 1973

By Encyclopædia Britannica, Inc.

Fifteenth Edition

© 1974, 1975, 1976, 1977, 1978, 1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986,
1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1997, 1998, 2002, 2003

By Encyclopædia Britannica, Inc.

© 2003

By Encyclopædia Britannica, Inc.

Copyright under International Copyright Union
All rights reserved under Pan American, Berne
and Universal Copyright Conventions by
Encyclopædia Britannica, Inc.

No part of this work may be reproduced or utilized
in any form or by any means, electronic or mechanical,
including photocopying, recording, or by any
information storage and retrieval system, without
permission in writing from the publisher.

Printed in U.S.A.

Library of Congress Control Number: 2002113989
International Standard Book Number: 0-85229-961-3

Britannica may be accessed at <http://www.britannica.com> on the Internet.

CONTENTS

1	LIGHT
29	LIMA
33	LINCOLN
40	LINGUISTICS
72	LISBON
77	The Art of LITERATURE
216	The History of Western LITERATURE
221	LOCKE
225	The History and Kinds of LOGIC
283	LONDON
299	LOS ANGELES
305	LUTHER
314	LUXEMBOURG
319	MADAGASCAR
329	MADRID
334	MALTA
339	MAMMALS
460	MANCHESTER
464	MANILA
468	MAO ZEDONG
473	MAPPING AND SURVEYING
495	MARKETING AND MERCHANDISING
509	MARKETS
526	MARSEILLE
531	Marx and MARXISM
544	MASKS
552	MATERIALS SCIENCE
566	The Foundations of MATHEMATICS
575	The History of MATHEMATICS
612	MATTER: Its Properties, States, Varieties, and Behaviour
691	MAXWELL
693	MEASUREMENT SYSTEMS
698	MECCA AND MEDINA
702	MECHANICS: Energy, Forces, and Their Effects
774	MEDICINE
829	MELBOURNE
832	MEMORY
841	MENTAL DISORDERS and Their Treatment
860	The History of Ancient MESOPOTAMIA
893	METABOLISM

The first part of the document discusses the importance of maintaining accurate records of all transactions. It emphasizes that every entry should be supported by a valid receipt or invoice. This ensures transparency and allows for easy verification of the data.

In the second section, the author outlines the various methods used to collect and analyze the data. This includes both primary and secondary data collection techniques. The primary data was gathered through direct observation and interviews, while secondary data was obtained from existing reports and databases.

The third section details the statistical analysis performed on the collected data. This involves the use of descriptive statistics to summarize the data and inferential statistics to test hypotheses. The results of these analyses are presented in a clear and concise manner, highlighting the key findings of the study.

Finally, the document concludes with a discussion of the implications of the findings and suggestions for future research. It notes that while the current study provides valuable insights, there are still several areas that require further investigation. The author hopes that this work will serve as a foundation for other researchers in the field.

Light

Light, a basic aspect of the human environment, cannot be defined in terms of anything simpler or more directly appreciated by the senses than itself. Light, certainly, is responsible for the sensation of sight. Light is propagated with a speed that is high but not infinitely high. Physicists are acquainted with two methods of propagation from one place to another, as (1) particles and as (2) waves, and for a long time they have sought to define light in terms of either particles or waves. In the early 19th century a wave description was favoured, though it was difficult to understand what kind of wave could possibly be propagated across the near-vacuum of interstellar space and with the extremely high speed of 300,000 kilometres per second (186,000 miles per second). In the latter half of the 19th century a British physicist, James Clerk Maxwell, showed that certain electromagnetic effects could be propagated through a vacuum with a speed equal to the measured speed of light. Thus, in the second half of the 19th century, light was described as electromagnetic waves (see ELECTROMAGNETIC RADIATION). Such waves were visualized as analogous to those on the surface of water (transverse waves) but with an extremely short wavelength of about 500 nanometres (one nanometre is 10^{-9} metre). The analogy is valid up to a certain point but the experimental results obtained at the end of the 19th century and in the early years of the 20th century revealed properties of light that could not have been predicted from knowledge that was obtainable about other waves. These results led to the quantum theory of light, which in its primitive form asserted that, at least in regard to its emission and absorption by matter, light behaves like particles rather than waves. The results of certain important experiments on the spreading of light into shadows and other experiments (on the interaction of beams of light) that supported the wave theory found no place in a particle theory. For a time it was believed that light could not be adequately described by analogy with either waves or particles—that it could be defined only by a description of its properties. A reconciliation of wave and particle concepts did not emerge until after 1924.

Two properties of light are, perhaps, more basic and fundamental than any others. The first of these is that light is a form of energy conveyed through empty space at high velocity (in contrast, many forms of energy, such as the chemical energy stored in coal or oil, can be transferred from one place to another only by transporting the matter in which the energy is stored). The unique property of light is, thus, that energy in the form of light is always moving,

and its movement is only in an indirect way affected by motion of the matter through which it is moving. (When light energy ceases to move, because it has been absorbed by matter, it is no longer light.)

The second fundamental property is that a beam of light can convey information from one place to another. This information concerns both the source of light and also any objects that have partly absorbed or reflected or refracted the light before it reaches the observer. More information reaches the human brain through the eyes than through any other sense organ. Even so, the visual system extracts only a minute fraction of the information that is imprinted on the light that enters the eye. Optical instruments extract much more information from the visual scene; spectroscopic instruments, for example, reveal far more about a source of light than the eye can discover by noting its colour, and telescopes and microscopes extract scientific information from the environment. Modern optical instruments produce, indeed, so much information that automatic methods of recording and analysis are needed to enable the brain to comprehend it.

From the standpoint of wave motion, blue light has a somewhat higher frequency and shorter wavelength than red. In the quantum theory, blue light consists of higher energy quanta than the red.

The subject of light is so wide and its associations are so numerous that it cannot be accommodated within one article of reasonable length. There are three main divisions of the subject of light: physical optics, physiological optics, and optical instrumentation. This article deals primarily with physical optics, treating the nature and behaviour of light. It also discusses the interaction of light with matter and describes such phenomena as luminescence in considerable detail. Although electromagnetic theory is considered here, further elucidation may be obtained in the article ELECTROMAGNETIC RADIATION. The article SENSORY RECEPTION includes the physiological and psychological aspects of light, while the section *Optical Engineering* in the article OPTICS treats the theory and technology of lenses, mirrors, and optical systems. The experimental evidence that led to the quantum theory of radiation is included in the present article along with a brief statement of some of the basic ideas. The quantum theory of radiation, however, is so closely associated with the quantum theory of matter that the two must be considered together, as is done in the section on *Quantum mechanics* of the MECHANICS article. (R.W.Di./Ed.)

The article is divided into the following sections:

General considerations	1	Dispersion and scattering	17
Historical survey	1	Mechanical effects of light	18
Basic concepts of wave theory	3	Quantum theory of light	19
Light spectrum	6	Photons	19
Velocity of light	7	The wave-particle nature of light	21
Interference and diffraction phenomena	8	Luminescence	22
Interference	8	Sources and process	22
Diffraction	11	Early investigations	22
Polarization and electromagnetic theory	13	Phosphorescence and fluorescence	23
Polarized light	13	Luminescence excitation	23
Electromagnetic-wave character of light	15	Luminescent materials and phosphor chemistry	24
The interaction of light with matter	16	Luminescence physics	25
Reflection and refraction	16		

General considerations

HISTORICAL SURVEY

From 500 BC to AD 1650. In this period, there were innumerable confusions and false starts toward an understanding of light. Sometimes an idea was stated, though

not clearly, and then almost forgotten for centuries before it reappeared and was generally accepted. The uses of plane and curved mirrors and of convex and concave lenses were discovered independently in China and in Greece. References to burning mirrors go back almost to the start of history, and it is possible that Chinese

Pythagorean hypothesis of light

and Greek knowledge were both derived from a common source in Mesopotamia, India, or Egypt. The formulation of general empirical laws and of speculation about the theory of light derives mainly from Mediterranean (Greek and Arab) sources. Pythagoras, Greek philosopher and mathematician (6th century BC), suggested that light consists of rays that, acting like feelers, travel in straight lines from the eye to the object and that the sensation of sight is obtained when these rays touch the object. In this way, the more mysterious sense of sight is explained in terms of the intuitively accepted sense of touch. It is only necessary to reverse the direction of these rays to obtain the basic scheme of modern geometrical optics. The Greek mathematician Euclid (300 BC), who accepted the Pythagorean idea, knew that the angle of reflected light rays from a mirror equals the angle of incident light rays from the object to the mirror. The idea that light is emitted by a source and reflected by an object and then enters the eye to produce the sensation of sight was known to Epicurus, another Greek philosopher of Samos (300 BC). The Pythagorean hypothesis was eventually abandoned and the concept of rays travelling from the object to the eye was finally accepted about AD 1000 under the influence of an Arabian mathematician and physicist named Alhazen.

Angles of incidence and of refraction—*i.e.*, the change in direction of a light ray going from one transparent medium to another—were measured by an astronomer, Ptolemy, in the 1st century in Alexandria. He correctly deduced that the ray is bent toward the normal (*i.e.*, the direction perpendicular to a boundary plane, such as the plane separating air and water) on entering the denser medium. A Dutchman, Willebrord van Roijen Snell, discovered the so-called sine law that gives the index of refraction (a measure of the change in direction) for light in a transparent medium. The laws of reflection and refraction were brought together by a 17th-century French mathematician, Pierre de Fermat, who postulated that the rays of light take paths that require a minimum time. He assumed that the velocity of light in a more dense medium is less than that in a less dense one in the inverse ratio of the indices of refraction.

Fermat's principle

The idea of rectilinear propagation of light—that is, that it travels in a straight line—was applied in a practical sense to drawing and painting long ago. Euclid was familiar with the basic idea, but the main theory was developed by Leonardo da Vinci, and a complete description of shadows was given by the Danish astronomer Johannes Kepler in 1604. Kepler also was the first to apply the laws of rectilinear propagation to photometry (the measurement of light intensities).

From 1650 to 1895. At the beginning of this period, the result of the conflict between the corpuscular theory and the wave theory was in doubt. At the end of the period, the wave theory was generally accepted and seemed capable of explaining all known optical phenomena though, with hindsight, it can now be seen that there were some important difficulties.

Diffraction—*i.e.*, the spreading of light into shadows—was first observed in Italy in the 17th century. In England, a worker, who independently noticed diffraction, also observed the interference colours of thin films, which are commonly seen today in an oil film on a wet road surface or in the iridescent colours of a butterfly's wing. He believed that light consists of vibrations propagated at great speed. Christiaan Huygens, of Holland, greatly improved the wave theory. In England, Sir Isaac Newton did not attach much importance to the small amount of spreading of light, and he knew that strictly rectilinear propagation could not be reconciled with the wave theory. Polarization phenomena (which can be accounted for by transverse wave motion in a single plane) discovered in the 17th century by a Danish physicist, Erasmus Bartholin, and by Huygens were not consistent with the theory of longitudinal waves (waves vibrating in the direction of propagation, like compression waves in a coiled spring), which was the only wave theory then considered. Newton therefore supported the corpuscular theory, although he did not reject the wave theory completely. He accepted a concept of a

luminiferous ether, and he postulated that the particles had "fits of easy reflection" and "fits of easy transmission"; *i.e.*, he assumed that they changed regularly between (1) a state in which they were reflected at a glass surface and (2) a state in which they were transmitted. He thus introduced periodicity—one of the basic ideas of wave theory—in a form that anticipates the quantum mechanics. Newton, using a glass wedge, or prism, discovered that white light can be separated into light of different colours and took the first steps toward a theory of colour vision.

In the century following his death the great authority of Newton was quoted to uphold the corpuscular theory and to oppose the wave theory in a way that he probably would not have approved. It was not until the 19th century that the work of Thomas Young of England; Augustin-Jean Fresnel, François Arago, and Armand-Hippolyte-Louis Fizeau, all of France; Irish scientist Humphrey Lloyd; and German physicist Gustav Kirchhoff established the transverse-wave concept of light; *i.e.*, light is a wave vibration at right angles to the direction of travel. A universal medium pervading all space and called the ether was supposed to be some kind of elastic solid. This made it possible to accept the transmission of light through a vacuum, but there was no completely satisfactory theory of the ether or of the way in which light is modified by transparent materials like glass. The necessity for an elastic solid disappeared when Maxwell proposed an electromagnetic theory of light. He stated the laws of electromagnetism in a clear mathematical form and generalized the concept of an electric current. From his equations he predicted the existence of transverse electromagnetic waves having a constant speed *c in vacuo*. The constant *c* had a value of 300,000 kilometres per second and was derived from measurements on electrical circuits. It was known from the work of Ole Rømer, a Danish astronomer; Jean-Bernard-Leon Foucault of France; and others that the velocity of light was not much different from the velocity constant *c*. A.A. Michelson, a physicist in the United States, measured the velocity of light and showed that it is equal to *c* within a small margin of experimental error. This result, together with the work of a German physicist (Heinrich Rudolf Hertz) on electromagnetic waves of larger wavelength, confirmed Maxwell's predictions (see ELECTROMAGNETIC RADIATION). The existence of a connection between electromagnetism and light had, indeed, been demonstrated in England much earlier in the century by Michael Faraday, who observed the rotation of the plane of polarization of a beam of light by a magnetic field (Faraday effect).

From 1900 to the present. Maxwell's theory is a theory of waves in a continuous (*i.e.*, infinitely divisible) medium. The energy of the waves is also infinitely divisible so that an indefinitely small amount can be emitted or absorbed by matter. Classical physical theories of the 19th century had predicted that in such a system the energy in equilibrium would be distributed so as to give an equal amount to each mode (frequency) of vibration. Because a continuous medium has an infinite number of modes of vibration, and the atoms (which constitute matter) have only a finite number, all the energy of the universe would be transformed into waves of high frequency. Maxwell understood this difficulty, which was later most clearly stated in the Rayleigh-Jeans law (after two English physicists, Lord Rayleigh and Sir James Hopwood Jeans) of the radiation of a blackbody (a body in which the intake and output of energy are in equilibrium). The German physicist Max Planck demonstrated that it is necessary to postulate that radiant-heat energy is emitted only in finite amounts, which are now called quanta. At first, it was hoped to retain, without modification, the theory of light as electromagnetic waves in free space and to use the quantum concept only in relation to the interaction between radiation and matter. In 1905, however, Einstein showed that, in the photoelectric effect, light behaves as if all the energy were concentrated in quanta—*i.e.*, particles of energy now called photons. In the same year, Einstein published the theory of relativity, which modified the whole of physics and gave a special role to the velocity constant *c*. Because light, in some situations, behaves like waves and, in others, like particles, it is necessary to have

Transverse wave concept established

The role of the velocity of light in relativity theory

a theory that predicts when and to what extent each kind of behaviour is manifested. The main development of the quantum mechanics, which does precisely this, took place between 1925 and 1935.

Light from ordinary sources is emitted by atoms the phases of which are not correlated with one another, so that there is a random irregularity or incoherence between the waves emitted from different atoms. This places severe restrictions on the conditions under which the periodicity associated with wave theory can be observed. In England, Lord Rayleigh appreciated this effect and knew that, by the use of pinholes or slits and light of a narrow range of wavelength, effectively coherent light could be produced. For a long time, interest in this topic lapsed. About 1935 Frits Zernike, a Dutch physicist, and others extended the theory of coherence to include the concept of partial coherence. This appeared to be of practical importance only in a few rather special applications (e.g., in the Michelson stellar interferometer; see below *Interference*). A theory of stimulated emission, attributable to the work of Einstein and an English physicist, Paul A.M. Dirac, postulated that under certain conditions atoms could be made to radiate in phase so that highly coherent radiation could be maintained indefinitely. The practical realization of these conditions, previously thought to be impossible, was achieved in 1960.

A second major development in the theory of light in this century is the application of so-called Fourier transform methods (a mathematical treatment of light waves) to a wide range of optical problems and, especially, to the transfer of information in optical systems (see OPTICS).

Today, the theory of light has again reached a point at which all known terrestrial phenomena are included in one logical theory. The known unsolved problems concern the transmission of light over the vast distances of intergalactic space. Here the theory of light impinges on the science of cosmology.

BASIC CONCEPTS OF WAVE THEORY

In this section on the wave theory of light, those properties of light that are consistent with a wave theory are described using a minimum of mathematical formulation. It is convenient to introduce the basic concepts of wave theory in relation to mechanical systems. Below, in the section on *Interference*, and beyond, it will be necessary to consider results obtained by more sophisticated mathematical methods, such as Fourier analysis.

General characteristics of waves. *Periodicity in time and space.* If one end of a stretched rope is vibrated, a wave will run along the rope. Figure 1 (top) represents a profile of the wave—i.e., a “snapshot” of the displacement of the rope from its normal position. It gives the variation of this displacement (indicated by ξ) at different points

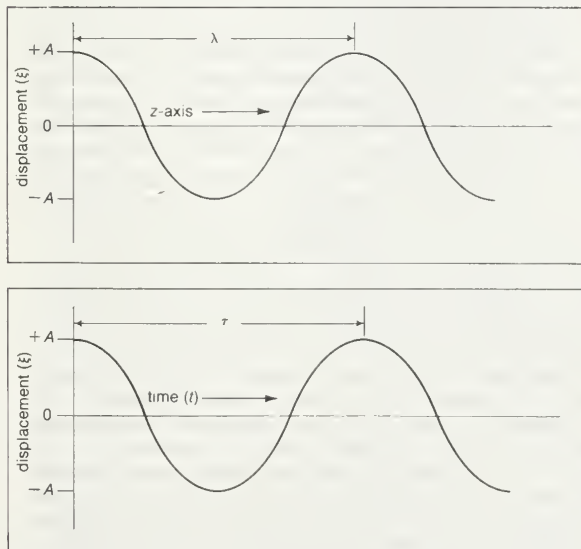


Figure 1: Wave profiles. (Top) Variation with position at one time. (Bottom) Variation with time at one place (see text).

(z) along the axis of propagation for one specific instant of time. Similarly, Figure 1 (bottom) shows the variation with time of the displacement at one arbitrary point on the axis. In Figure 1 (top) the distance between successive crests is constant and is called the wavelength (λ). Similarly, the constant time between crests in Figure 1 (bottom) is called the period (τ). The temporal frequency ($\nu_t = 1/\tau$) is the number of vibrations per unit time and the spatial frequency or wave number ($\nu_s = 1/\lambda$) is the number of waves per unit length. The wave shown in Figure 1 (top) may be represented by the cosine of an angle (ϕ) to give the displacement for a particular point on the axis at any instant of time:

$$\xi = A \cos \Phi = A \cos 2\pi(\nu_t t - \nu_s z), \quad (1)$$

in which ξ is the displacement at any point z on the axis at a time t , A is the amplitude (the maximum displacement); the angle ϕ (phi) in this case is equal to $2\pi(\nu_t t - \nu_s z)$ and is called the phase angle, or simply, the phase.

Energy. The energy per unit volume (W) stored in a wave motion is proportional to the square of the amplitude (A) so that, with a suitable choice of units, $W = A^2$.

Phase velocity. Any one crest moves forward a distance λ in a time τ ; i.e., with a velocity b of the wavelength divided by the period or the temporal frequency divided by the spatial frequency,

$$b = \frac{\lambda}{\tau} = \frac{\nu_t}{\nu_s} = \lambda \nu_t. \quad (2)$$

The velocity b is called the phase velocity because the phase angle ϕ will remain constant when the time t changes by an incremental amount t_0 and z changes by $z_0 = bt_0$. (This may be seen by substituting $t = t_0$ and $z = z_0$ in the expression for this phase and using $b = \nu_t/\nu_s$.)

The velocity of light in vacuum (denoted by c) is the same for all frequencies; all colours travel through space with the same speed. The phase velocity (denoted by b) in a material medium, on the other hand, depends on the medium and on the temporal frequency and, hence from equation (2), on the wavelength.

Wave surfaces. Two-dimensional waves are formed by vibrating (dipping) the end of a rod up and down in the surface of a liquid. Waves spread from the point of origin (where the rod contacts the surface) and, at any moment, the phase at any point on a circle is the same; i.e., if, at a given moment, the wave is at a maximum at one point on a circle then it is at a maximum everywhere on this circle, and the circle as a whole is a wave crest. Similarly, a trough is found at all points on another circle (the radius of which is $\lambda/2$ greater than that of the first circle). As the waves progress farther and farther from the origin, they become less strongly curved about the origin so that, at great distances, they are approximately plane waves.

Light waves are propagated in three dimensions and, for waves from a point source in an isotropic medium (i.e., one in which the speed is the same along any radius), the phase is constant over spherical surfaces drawn about the point source as a centre. The surfaces of constant phase are called wave surfaces, and waves are called plane, spherical, ellipsoidal, and so on according to the shapes of the wave surfaces.

Reflection and refraction. The similarity between the behaviour of light waves and the surface waves of a liquid may be demonstrated with the so-called ripple tank. For reflection of a train of surface waves incident on a flat object, it may be readily observed that the angle of reflection is equal to the angle of incidence. For waves that are refracted in passing from one medium of the ripple tank in which the phase velocity is b_1 , to another in which the phase velocity is b_2 , measurements of angles of incidence (θ_i) and refraction (θ_r) of the surface waves verify Snell's sine law of refraction; i.e., that the ratio of the sines of the angle of incidence and refraction is a constant, or

$$\frac{\sin \theta_i}{\sin \theta_r} = \frac{b_1}{b_2} = n_{12}, \quad (3)$$

in which the constant n_{12} is called the index of refraction from medium 1 to medium 2. The index of refraction (n)

from vacuum to a material medium is called the index of the medium and, for transparent mediums is always greater than unity (one). When n_{12} is less than unity, as happens when light is refracted as it passes from glass into air, the refracted ray grazes the surface if $\sin \theta_c = n_{12}$, θ_c being the angle of incidence in the glass. At angles of incidence greater than this critical angle there is total reflection; *i.e.*, light, instead of penetrating into the air, is reflected back into the glass.

Dispersion. Newton found that, when a beam of white light is refracted by a glass prism, it is dispersed, or split, into beams of different colours. This phenomenon is now interpreted in the following way: the velocity of light in glass varies fairly rapidly with its wavelength, whereas its velocity in air varies little; thus the index of refraction and hence the angle of refraction depend on wavelength. A beam of white light, containing as it does a wide range of wavelengths, is thus dispersed by a glass prism so that light of one wavelength emerges from it in a different direction from light of another wavelength. Because colour depends on wavelength, the emergent light forms a spectrum (see Plate). All material mediums are, to some extent, dispersive (*i.e.*, phase velocity varies with the temporal or spatial frequency).

Wave groups. When a stone is dropped into a quiescent pond, a few waves may be seen travelling out from the point of impact. This group of waves maintains its identity as it is propagated over a considerable distance, although it finally dies away. The velocity of the group as a whole is called the group velocity. Careful observation shows that the group velocity is less than the phase velocity. Individual waves may be seen to appear at the back of the group, advance through it, and die out as they reach the front of the group. In a nondispersive medium the group velocity is equal to the phase velocity, while in a dispersive medium it may be greater than, less than, or equal. For light waves, the group velocity is almost always less than the phase velocity.

Interference. When two or more wave motions are present at the same place and time, the simplest assumption is that the resultant displacement (ξ_R) is the algebraic sum of the individual displacements ($\xi_1, \xi_2, \xi_3, \dots$), *i.e.*,

$$\xi_R = \xi_1 + \xi_2 + \xi_3 + \dots + \xi_N. \quad (4)$$

Nearly all observations on light are in accord with this equation, which is a statement of the principle of superposition. These phenomena constitute the subject of what is known as linear optics. The possibility that additional phenomena might be observed at high intensities of light has long been accepted, and the use of lasers in the attainment of the necessary high intensities has led to the discovery of frequency doubling and other effects that cannot be predicted from equation (4). These new observations constitute the material of nonlinear optics (see OPTICS). Equation (4) is valid for all the phenomena of interference, diffraction, etc., which will be described in this article.

Two waves are said to be coherent if their phase difference remains constant during a period of observation. Figure 2 shows two equal coherent plane waves travelling

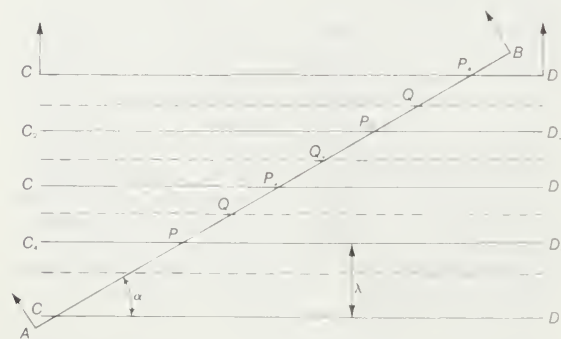


Figure 2: Interference of two plane waves AB and CD with directions inclined at an angle α . The crests of CD are represented as C_1D_1, C_2D_2, \dots , and the troughs are shown as broken lines (see text).

across the same space, with the wave fronts inclined at a small angle α , AB representing a surface corresponding to a crest of one wave. (The surface must be assumed to be perpendicular to the page.) C_1D_1, C_2D_2, \dots , represent surfaces that correspond to crests of the other wave. The intermediate dotted lines represent troughs. At points such as P_1 (and P_2, P_3, \dots), a crest of one wave coincides with a crest of the other and according to the principle of superposition the displacement is twice that of either wave alone. At points Q_1, Q_2, \dots , a crest of one wave meets a trough of another; so the displacements being equal and opposite, the resultant is zero. Thus, an observer looking at a plane that is perpendicular to the page and passes through AB sees a series of straight lines through P_1, P_2, P_3, \dots , etc., representing large displacement and a series of lines through Q_1, Q_2, Q_3, \dots , etc., representing zero displacement.

There are many ways in which coherent beams of light can be made to cross at an angle of about one part in a thousand. The eye (or a low-power magnifier) can be focussed on a plane such as that through AB . The resulting parallel light and dark lines are called interference fringes (Figure 3). From Figure 2 it may be seen that the separation (d) of two bright fringes is λ/a or $1,000 \lambda$ if $a = 0.001$. When a has this value, $d = 0.5$ millimetre for blue-green

Interference fringes

Milward T. Rodine

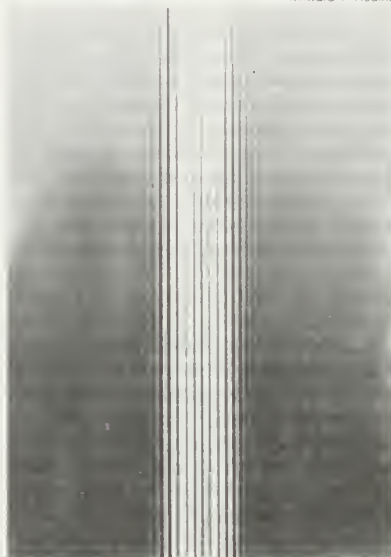


Figure 3: Two-beam interference fringes from Young's double slits or Fresnel's biprism (see text).

light and this would imply that λ is about 0.5×0.001 or $1/2,000$ part of a millimetre (this is usually written 500 nanometres).

In this experiment the spatial periodicity of the light waves (about 2,000 waves per millimetre) has been made to produce fringes with periodicity of about two per millimetre. The spatial periodicity of a light wave is too high for the human eye, and it cannot be magnified directly. Interference methods effectively magnify it so that the resultant fringes can be seen by eye or with a convenient magnification. The following method of producing interference fringes, developed by Thomas Young, is now called Young's experiment.

In the arrangement shown in Figure 4, light of one wavelength passes through a slit S producing semicylindrical waves that are intercepted by two other slits P_1 and P_2 . The two slits P_1 and P_2 act as secondary sources of coherent, semicylindrical waves the combined effect of which is observed on the plane perpendicular to the page and designated AB . In a typical case the separation (a) of P_1 and P_2 is a millimetre and the distances l_1 and l_2 are each about a metre. The slits are a centimetre or so long but are much less than a millimetre wide. They are accurately parallel to one another and, as represented in the drawing, are at right angles to the page. Because the waves from P_1 and P_2 are indirectly derived from the same small source, they are coherent. When they cross plane AB they are

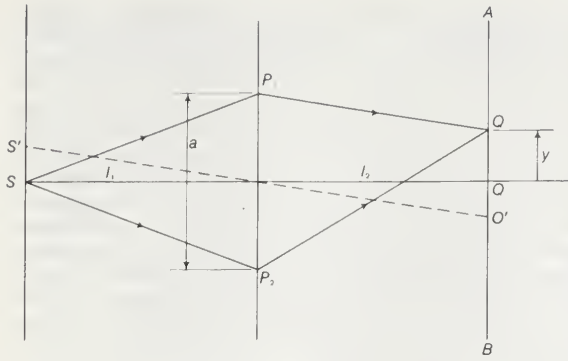


Figure 4: Young's experiment (see text).

nearly plane because of the large radius, and they intersect at an angle α equal to 0.001. It may be shown that the intensity (I) for these fringes varies from point to point along the line AB in the way shown in Figure 5 (curve A), which is in accord with the equation

$$I = 2A^2(1 + \cos \epsilon) = 2I_0(1 + \cos \epsilon), \quad (5)$$

in which A is the amplitude of either wave, I_0 is the intensity of one wave acting alone and the phase difference $\epsilon = 2\pi ya/\lambda l_2$. Bright fringes are seen in positions for which $\epsilon = 2\pi p$ or $y = p\lambda/l_2 a$ (in this case p is a whole number, which may be positive, zero or negative—0, ± 1 , ± 2 , ± 3 , etc.). Because $\cos \epsilon$ varies from -1 to $+1$, I varies from $4I_0$ to zero. The average, in accordance with the law of conservation of energy, is $2I$.

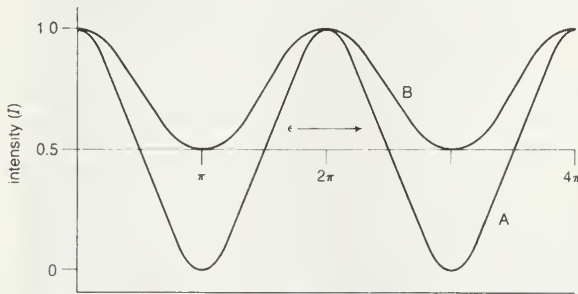


Figure 5: Interference fringes obtained in Young's experiment (see text).

Diffraction. Plane waves that pass through a restricted opening emerge as divergent waves. When the opening is less than one wavelength in diameter the emergent wave is nearly spherical. Whenever a beam of light is restricted by holes or slits or by opaque obstacles that block out part of the wave front, some spreading occurs at the edges of geometrical shadows. This effect, called diffraction, is also obtained with transparent obstacles that cause an irregularity in the wave front. Diffraction can be demonstrated by allowing a parallel beam of light to fall on a grating consisting of an array of equally spaced narrow slits. If the extent of physical separation of two adjacent slits is e , then the path difference between any two adjacent rays emitted in a direction symbolized by θ is $e \sin \theta$, and if this path difference is an integral number (p) of wavelengths,

$$e \sin \theta = p\lambda, \text{ or } v_s \sin \theta = pg, \quad (6)$$

in which v_s is the spatial frequency ($1/\lambda$) and g is the number of lines per unit width of the grating, then the waves from different slits have phases that differ by angles of $2p\pi$, and they reinforce one another. Thus, when lenses are employed with a grating, sharp lines are obtained for each wavelength at values of θ corresponding to integral values of p . If white light is used, each line is drawn out into a spectrum of wavelengths because the direction of reinforcement depends on the wavelength.

Polarization. In the propagation of waves on a rope or across the surface of a liquid the displacement (as shown in Figure 1) is in a direction perpendicular to the direction of propagation and the waves are said to be transverse. Sound waves in a gas consist of alternate dilation and compression and the displacement is in the direction of

propagation. The waves are longitudinal. If a beam of longitudinal waves is propagated in a vertical direction, there is nothing to distinguish one azimuthal plane from another—everything that is true for an east-west plane is equally true for a north-south plane. With transverse waves the displacement may be in the east-west plane; in that case, there is no component in the north-south plane, and this should manifest itself in the form of a property that depends on the azimuth. Such an effect is called an azimuthal property. An ordinary beam of light from a thermal source does not exhibit any azimuthal property, but experiments show that light can have an azimuthal property and must be represented by transverse waves.

Azimuthal property

If an unsilvered glass plate has an index of refraction equal to 1.5 and the angle of incidence of a beam of light is 57° , about 15 percent of the light will be reflected from the two glass surfaces of the plate (Figure 6); this percentage will not be altered when the glass plate is rotated about an axis parallel to the beam of light so as to change the azimuth of the plane of reflection. If a second mirror (G_2), parallel to the first (G_1), is used to reflect the beam in the same plane as that of the original reflection, about 30 percent of the light incident on the second plate of glass will be reflected; but if the second plate is turned so as to reflect the light in a plane perpendicular to that of the first reflection—*i.e.*, out of the plane of the page—hardly any light will be reflected. Thus, after the first reflection, the beam of light will have acquired an azimuthal property—it will be reflected more strongly when the transverse displacement is in one azimuthal plane than when in another. Further tests will show that the transmitted light has a complementary azimuthal property; it is more strongly reflected in the perpendicular plane—though the difference is less marked.

These results may be understood if ordinary light consists of a mixture of transverse waves with displacements in all azimuthal planes but only one component is reflected from a glass surface when the angle of incidence is 57° . The reflected light is said to be plane-polarized because all of the displacement of the wave is in one azimuthal plane. The transmitted light (about 85 percent of the whole) contains about 50 parts of a component that is polarized in a perpendicular plane and about 35 parts of light that is polarized in the same way as the reflected light. It is more strongly reflected in the plane of the page, but because it is only partially polarized, the azimuthal effect is less.

From R W Ditchburn, *Light* (1963), Interscience Publishing, Inc., by permission of John Wiley & Sons, Inc.

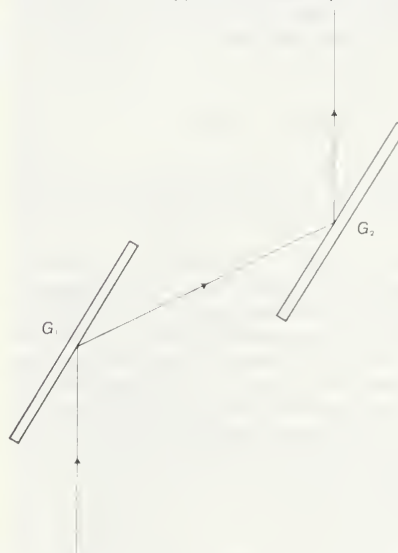


Figure 6: Malus' experiment. Successive reflections at two unsilvered mirror surfaces, G_1 and G_2 (see text).

The above experiments do not show whether or not the reflected light has its displacement in the plane of reflection or perpendicular to it. It is a matter of choice whether the reflected light is said to be polarized in or perpendicular to the plane of reflection. Some controversy (and some

difference of nomenclature) that formerly led to confusion was removed by the electromagnetic theory (see below). In this theory light is represented by two vectors (quantities that can be represented graphically by arrows that point in the field directions), a magnetic vector in the plane of reflection and an electric vector perpendicular to it. Confusion is avoided by specifying the plane of the electric vector instead of speaking of the plane of polarization.

The azimuthal property of reflected light at the surface of any medium—glass, plastic, a liquid—is most strongly manifested when the angle of incidence is so chosen that its tangent is equal to the index of refraction; that is, it satisfies Brewster's law (after Sir David Brewster, a British physicist), which states that, at the polarizing angle, the incident and refracted beams make an angle of 90° with one another; $\tan \theta_i = n$, in which θ_i is the angle of incidence, called the polarizing angle, and n is the index of refraction of the medium. Nevertheless, there is some azimuthal difference after reflection at any angle except $\theta_i = 0$ or $\theta_i = 90^\circ$. Other ways of producing polarized light are described in a later section.

It is found that the plane of polarization of a beam of polarized light is rotated when the beam is passed through certain mediums (especially sugar solutions). These mediums are said to be optically active. Most mediums do not normally rotate the plane of polarization, but do so when there is a magnetic field in the direction of propagation (the Faraday effect).

Optical activity

The wave equation. The expression for a plane wave, given in equation (1) and showing the relationship between displacement (ξ), the time span (t), and distance (z) along the wave, may be differentiated twice with respect to t and z ; that is, to find out how the displacement changes with position and time. This operation yields the partial differential equation:

$$\frac{\partial^2 \xi}{\partial z^2} = \frac{1}{b^2} \cdot \frac{\partial^2 \xi}{\partial t^2} \quad (7a)$$

in which b is the phase velocity. For a three-dimensional wave the analogous expression is

$$\frac{\partial^2 \xi}{\partial z^2} + \frac{\partial^2 \xi}{\partial y^2} + \frac{\partial^2 \xi}{\partial x^2} = \frac{1}{b^2} \cdot \frac{\partial^2 \xi}{\partial t^2} \quad (7b)$$

There are many solutions of this basic equation. Some correspond to the sinusoidal plane waves, which have already been considered. Others correspond to groups of plane waves that differ slightly either in direction, or wavelength, or both. Yet another solution of the general wave equation is:

$$\xi = \frac{A}{r} \cos 2\pi(\nu t - \nu_1 r), \quad (8)$$

in which r is the magnitude of a radius vector drawn from the origin and A is a constant. This represents spherical waves.

Energy of a beam of light. The energy in a small volume (dV), through which plane waves are passing, is proportional to the product of the square of the amplitude (A), or its energy per unit volume (W), and the small volume; that is, $A^2 dV = W dV$. The rate of transport of energy across a surface normal to the direction of propagation is proportional to the product of the energy per unit volume, the phase velocity, and a small area (dS) normal to the direction of propagation, or $W b dS$. For spherical waves, the rate of transport is inversely proportional to r^2 , i.e., $(A/r^2) dS$. Because the area of a sphere is $4\pi r^2$ in which r is its radius, this equation implies that the total energy crossing any sphere surrounding a point source is independent of the radius. Thus, inverse-square law for the intensity of radiation at a distance r from a point source is in accord with the law of conservation of energy—the total energy of a wave remains the same even though the wave is spread over a greater area.

Doppler-Fizeau effect. The length of a wave train emitted in one second by a stationary light source is equal to the velocity of light (c) times one second, which in

itself is equal to the product of its frequency (ν) times its wavelength (λ)—i.e., $c = \nu\lambda$. If the source moves away from the observer with a velocity (v) that is small compared with the velocity of light, then the length of the wave train increases so as to be numerically equal to the sum of the two velocities ($c + v$) and the number of waves remains the same. The wavelength λ increases to λ' by a factor $(c + v)/c$; that is $\lambda' = (1 + v/c)\lambda$. This change was discovered by an Austrian physicist, Christian Doppler, in the 19th century in relation to sound waves and subsequently applied to light waves by Fizeau. It is called the Doppler-Fizeau effect. The Doppler-Fizeau effect is easily observed when part of the light from a gas laser is allowed to be scattered by a moving body and mixed with a little unscattered light. It is known from the study of sound waves that the beat frequency is equal to the difference between the frequencies of the two waves that are mixed. Although the frequency of light waves is extremely high (more than 10^{14} per second), the beat frequency may be a megahertz (10^6 cycles per second), which is easily detected by radio amplifiers, or even a few hundred cycles per second, which the human ear can detect. Thus, just as interference fringes provide a periodic phenomenon in which two light waves combine to produce fringes of low spatial frequency, so the Doppler-Fizeau effect produces beats the temporal frequency of which is a known, but very small, fraction of the temporal frequency of the light waves. In this way the periodicity of light in both space and time is exhibited and measured.

Beat frequency

LIGHT SPECTRUM

It was seen, in the preceding section, that white light can be dispersed into a spectrum by refraction, by diffraction, or by interference. Newton showed that if a suitably oriented slit is used to select a small region of the spectrum, the light that passes through the slit is much more homogeneous than the original white light, and he was unable to observe any further dispersion when passing this light through a second prism. Delicate methods of interferometry nevertheless show that this light is never entirely of one wavelength, however fine the slit, but covers a range ($\Delta\lambda$) of wavelengths. The ratio of the wavelength divided by this range, which measures the purity of the spectrum, may be a few thousand for a spectrum formed by a prism and up to a million for a spectrum formed by a large diffraction grating. It is never infinite, as it would be if $\Delta\lambda$ were zero.

The spectrum of a hot body such as the solar photosphere is continuous (every wavelength is represented); but a German physicist, Joseph von Fraunhofer, early in the 19th century observed that the solar spectrum contains numerous dark lines appearing at certain wavelengths, which are attributed to wavelengths originally emitted by inner layers of the Sun but then absorbed by various elements (in gaseous form) in the cooler outer layers (see Plate). Emission spectra produced by electric sparks and arcs contain sharp bright lines which are characteristic of the elements in the electrodes.

In monochromatic light, colour and wavelength are associated. Nevertheless, as Newton said, "the rays, to speak properly, are not coloured." Colour is a sensation in the human mind. Light of one wavelength can stimulate the visual system so that a certain colour sensation (e.g., red) is produced. The way in which the visual system analyzes colour is entirely different from the way in which physical instruments form a spectrum (see SENSORY RECEPTION).

There are a number of ways in which spectra are produced in nature. The rainbow is the most striking of these. The primary rainbow is formed by reflection and refraction of light in raindrops. The rays emerging from the drops are spread out, but for any given wavelength there is a minimum angle of deviation and there is a concentration of energy at this angle. For green light the minimum angle of deviation is about 138° and an observer with his back to the Sun sees the bow at an angle of 42° to the direction of the Sun's rays. Because of the dispersion of water, the angles for different wavelengths are not exactly the same, and the red is seen on the outside and blue on the inside of the bow. A weaker rainbow is formed by

Primary and secondary rainbows

rays that have been twice reflected. In this the colours are reversed. Still weaker supernumerary bows are caused by diffraction in droplets. A rainbow may be regarded as a spectrum of the Sun, but the purity is low.

VELOCITY OF LIGHT

The accepted value of the velocity of light (c) in vacuum is 299,792.458 kilometres per second (see Table 1). The velocity is the same for all wavelengths over the whole range of the electromagnetic spectrum from radio waves to gamma rays. Methods of measurement are of three types: (1) measurement of the time (T) in which a group of waves covers a known distance (l), (2) measurement of the frequency (ν) and wavelength (λ) of monochromatic waves, and (3) indirect methods, such as measurement of the change of frequency or wavelength (Doppler-Fizeau effect) when a beam of light is reflected from a mirror moving with a known velocity.

Table 1: The Constant c (in kilometres per second)		
	year	value
Derived from measurements of the velocity of light		
Michelson	1927	299,796 \pm 4
Michelson, Pearson and Pease	1935	299,774 \pm 11
Value accepted in 1941	1941	299,773 \pm 3
Bergstrand	1951	299,793.1 \pm 0.2
Bergstrand (mean value)	1957	299,792.9 \pm 0.2
Value adopted by 17th General Congress on Weights and Measures	1983	299,792.458
Derived from measurements on radio waves		
Essen (10 ⁴ MHz)	1950	299,792.5 \pm 1
Froome (2.4 and 7.5 \times 10 ⁴ MHz)	1951-58	299,792.5 \pm 0.1
Value adopted by 12th General Assembly of the Radio-Scientific Union	1957	299,792.5 \pm 0.4
Derived from electrical measurements		
Rosa and Dorsey (ratio of units)	1907	299,788 \pm 30
Mercier (Lecher wires)	1923	299,795 \pm 30

Methods of type (3) have, so far, given an accuracy of only a few percent. Methods of type (2) cannot be used for light waves because the frequency is about 1.5×10^{14} hertz and is too high to be measured directly. The remainder of this section will review measurements of the velocity of light by methods of type (1) and compare the results of the best measurements with the results obtained for radio waves by methods (1) and (2).

Astronomical measurements. In 1676 Rømer made careful measurements of the times at which satellites of Jupiter were eclipsed by the planet. The times observed did not agree with those calculated on the assumptions of a constant period of rotation and of instantaneous transmission of light. Starting at a time when the Earth was at its nearest to Jupiter, the apparent period increased and the eclipses became increasingly later than the calculated times as the Earth receded from Jupiter. Similarly, the period shortened when the Earth was moving toward Jupiter. The observed times were consistent with a finite velocity of light such that the time for it to transverse the Earth's orbit is about 1,000 seconds. Taken with modern values of the size of the Earth's orbit, the derived value of the velocity is 298,000 kilometres per second. It is remarkable that this first measurement was even of the correct order; the most important conclusion was that the velocity of light is finite. An English astronomer, James Bradley (died 1762), obtained a similar value by the so-called aberration method, based on the apparent motion of stars as the Earth travels in its orbit about the Sun.

Early terrestrial experiments. In terrestrial experiments by method (1), the beam of light is periodically marked either by interrupting it at regular intervals or by modulating it (alternately increasing and decreasing its intensity). The marked beam is transmitted to a distant mirror and the return beam passes through the apparatus that interrupts or modulates the outgoing beam and then to a detector. If the time required for transmission to the distant mirror

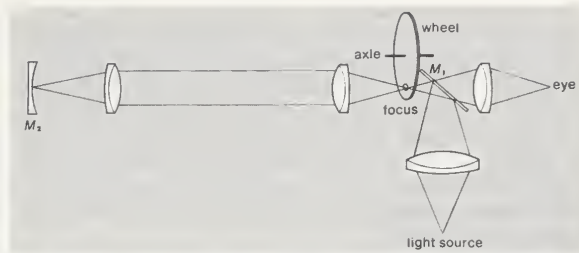


Figure 7: Fizeau's method for measuring the velocity of light.

From R.W. Ditchburn, *Light* (1963), Interscience Publishing, Inc., by permission of John Wiley & Sons, Inc.

and return is $\frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \dots$ times the period of the interrupter (or modulator), then the amount that reaches the detector is small. It is usual to adjust either the path length or the period of the interrupter or modulator until the light registered by the detector is a minimum. In the earlier experiments, a mechanical chopper was used as interrupter, and the eye was the detector. Later experimenters used electronic modulators and photoelectric detectors.

The apparatus used by Fizeau in 1849 is shown in Figure 7, in which M_1 is a partially reflecting mirror and M_2 is a fully silvered mirror. As the speed of the wheel (which has 720 teeth) was increased from zero, it was found that the light was first eclipsed by a tooth when the speed was about 12.6 revolutions per second—*i.e.*, when the time to make the round trip was 560 microseconds (0.00056 second), the length of the double path being 17.3 kilometres (about 10 miles). The chief error in the measurement lay in the difficulty of determining the exact speeds at which the light received by the eye at E was at a minimum. Essentially the same method was used by others between 1874 and 1903. The accuracy gradually improved, and it was shown that the velocity is between 299,000 and 301,000 kilometres per second.

In 1834 Sir Charles Wheatstone of England suggested a method incorporating a rotating mirror for interrupting the light that was later developed by Arago (1838) and Foucault (1850). It was considerably improved by Michelson, who made measurements from 1879 to 1935.

Use of a rotating mirror

Michelson's measurements. Figure 8 shows the arrangement used in 1927. The mirror M_3 is a little above the plane of the diagram, and M_3' is a little below. Light from the source S passes to one face of the octagonal mirror M_1 and then to M_2 , M_3 , and M_4 . From M_4 it goes to the mirror M_5 at a distance of about 35 kilometres (about 22 miles). It returns via M_6 , M_4 , M_3' , and M_1 to the octagon. An image of S is seen in an eyepiece at E . The octagonal mirror rotated at 528 revolutions per second. It turned through approximately one-eighth of a revolution during the transit of the light. If the rotation were exactly one-eighth of a revolution, the image would be undisplaced from the position it had when the mirrors were stationary. In some of Michelson's experiments, the speed of rotation was slowly changed until this condition was obtained. In others, the speed and distance were fixed, and a small displacement of the image was measured.

It is difficult to estimate the accuracy of Michelson's 1927 and 1935 experiments, and it is no longer important to do so in view of the more accurate measurements made since 1945. His most important contribution to the measurement of the velocity was the proof that the velocity agreed with Maxwell's prediction to better than one part in a thousand. This gave confidence to those working on applications of the electromagnetic theory.

From R.W. Ditchburn, *Light* (1963), Interscience Publishing, Inc., by permission of John Wiley & Sons, Inc.

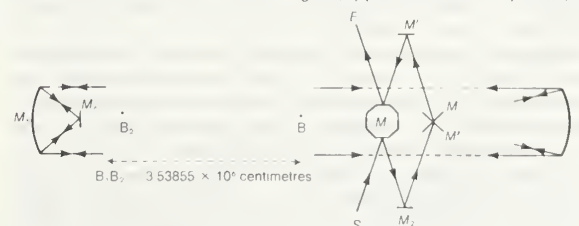


Figure 8: Michelson's Mount Wilson experiment, 1927.

The electro-optical shutter. This device, based on the Kerr effect (see below), makes it possible to modulate a beam of light at frequencies more than 10,000 times the highest frequency of interruption used by Michelson and obtain values in reasonably good agreement with each other and with Michelson's later work. This method was greatly improved by E. Bergstrand in Sweden, who reduced the random errors by a factor of more than 30 and obtained a value for the velocity of light of 299,793.1 kilometres per second.

Radio-frequency measurements. The velocity of electromagnetic waves of radio frequency in vacuum has been measured by several methods. An English physicist, Louis Essen, measured (1950) the resonance frequency of a cavity resonator (an electromagnetic device) whose dimensions were also determined with high accuracy. Keith Davy Froome, a physicist in England, measured (1952 and 1958) the wavelength in air, corresponding to a known frequency, using a microwave interferometer. The results of these and other measurements are in agreement with those of Bergstrand to within a few parts per million. The velocity of radio waves in vacuum is thus equal, within this accuracy, to the velocity of light. The velocity of gamma rays is also the same, within the much lower accuracy of this last measurement. Table 1 summarizes the measurements of the velocity constant (c) and shows that there is now satisfactory agreement between results obtained over a wide range of conditions.

Since the publication of the special theory of relativity (1905), the constant c has been recognized as one of the fundamental constants of modern physics. For this reason, attempts will undoubtedly be made to measure it with even greater precision. The use of lasers may help, but a major improvement will require the establishment of better standards of length and time than those now available.

Velocity in material mediums. All measurements of the velocity of light involve interruption or modulation of a beam of light so as to form groups of waves and the velocity measured is the group velocity. The difference in magnitude between the wave velocity and the group velocity of light in air is only about one part in 50,000, but in most glasses and in some liquids it is much larger. Michelson obtained 1.758 for the ratio of the velocity in air to the velocity in carbon disulfide. The inverse ratio of their indices of refraction is 1.64 and the value calculated from this for the ratio of group velocities is 1.745 for wavelength 580 nanometres, close to Michelson's observations. Bergstrand found that the ratio of the velocity *in vacuo* to the velocity in a certain glass was 1.550 ± 0.003 . The refractive index of the glass was 1.519, but the ratio of c to the group velocity was 1.547. The experimental results thus agree with those calculated on the assumption that the measured velocity is the group velocity.

Interference and diffraction phenomena

INTERFERENCE

Quasi-monochromatic waves. A perfectly monochromatic wave, represented by equation (1), has constant amplitude and is not limited in space or in time. Sources of light (other than lasers) emit waves the amplitude of which varies with time. For example, a single undisturbed atom emits a damped wave (Figure 14A). Under favourable conditions the damping is so weak that 10^7 waves are emitted before the amplitude has fallen to half its initial value and the change of amplitude is not significant over a distance of several thousand wavelengths. Wave trains of this type are said to be quasi-monochromatic. Superposition of these waves gives interference when the path difference is not too large.

Photometric summation. Figure 5, curve A, shows the way in which the intensity of light varies from place to place when two monochromatic or quasi-monochromatic waves overlap. The intensity at a point in the region where the waves overlap may be expressed as the sum of two terms: (1) the sum of the intensities of each wave acting alone ($2I_0$ if each alone would give intensity I_0); (2) a term representing the interference of the waves. The second term varies from point to point along the direction

of propagation between the values $-2I_0$ and $+2I_0$. Thus the total intensity varies from $4I_0$ (i.e., twice the intensity sum) to zero. Now, when a large number of waves from different sources cross a certain space, the fringes caused by the interference of each pair of waves have their maxima in different places and the overall result is that, at any point, the interference terms are positive nearly as often as they are negative and their total sum is nearly zero. In this case the resultant intensity at any point caused by a number of sources is just equal to the sum of the intensities (at that point) of each source acting alone. This is the law of photometric summation and is used by illumination engineers in calculating the illumination on a surface that receives light from various sources. Interference fringes are obtained only when experimental conditions are such that the interference fringes caused by light emitted from different atoms all have their maxima in the same places (or near to the same places). The interference term then becomes a significant fraction of the summation term. This may be achieved either (1) by using two secondary sources (such as the two slits used in Young's interference experiment), which are both derived from the same primary source, or (2) by using a laser in which the source atoms are stimulated in such a way that the phase relations between them remain constant during the period of observation.

Visibility of interference fringes. The distribution of intensity in interference fringes, shown in Figure 5, curve A, represents an ideal that is closely approached in some experiments, but generally the distribution is such that the fluctuations that constitute the fringes are superposed upon a nearly uniform background. Michelson defined the visibility of fringes as the difference between the maximum and minimum intensity of a fringe divided by their sum, or

$$V = \frac{I_{\max} - I_{\min}}{I_{\max} + I_{\min}} \quad (9)$$

in which V is the visibility, I_{\max} is the maximum intensity, and I_{\min} is the minimum. The fringe visibility is thus always between zero and one. When the minimum intensity is zero, the visibility equals one. Obviously, fringes for which V is less than one are obtained when waves of unequal amplitude are superposed because the weaker cannot, at any point, annul the stronger. It is also found, however, that even when the intensities are equal, the visibility is usually less than one (as shown in Figure 5, curve B). Further consideration of Young's slit experiment leads to recognition of two conditions that must be fulfilled to obtain fringes of high visibility. These relate to their geometrical condition and spectral range.

Geometrical conditions. In the arrangement shown in Figure 4, the centre of the fringe system is at a position O on the screen, on the straight line from the source slit S to a position midway between P_1 and P_2 (the slits are all assumed to be extremely narrow). If slit S is moved to S' then the centre moves to O' . If, instead of moving the slit S to this new position, it is gradually widened, the intensity at any point Q is found by adding the intensities of waves emitted by atoms behind different parts of the slit. Because the fringes on plane $O'OQ$ produced by light from different parts of slit S are not in register, there cannot be zero intensity at any point in the pattern. As slit S is widened the fringes gradually become blurred—i.e., the visibility falls from unity to zero. If $l_1 = l_2$, no fringes are seen when the width of slit S is about equal to the distance (d) between successive fringes.

Spectral range. In the case in which the slit S is extremely narrow and the light is not all of exactly one wavelength, the path difference and the phase difference will be zero at the centre of the fringe system for all wavelengths, so that for all wavelengths there is maximum intensity at the centre O of the fringe system. Because the separation of the fringes is proportional to the wavelength, the fringes produced by light of different wavelengths gradually go out of register as the path difference is increased. With white light, one clear fringe is seen in the centre. A few coloured fringes are seen on either side because the

Recognition of the velocity of light as a fundamental constant

Two conditions for high visibility

eye makes a certain degree of separation of the colours. If a filter is used to restrict the light to a band of say 50 nanometres wide, then about ten fringes may be seen on either side, and this number is increased if the wavelength range is further restricted.

These two causes of reduced visibility differ in that the geometrical condition affects all parts of the fringe system equally and the effect of the spectral range increases as the path difference increases. In discussing these phenomena it has been assumed, in accordance with the preceding discussion, that the intensity of light from different atoms obeys the law of photometric summation. It is also assumed that the photometric law applies when different wavelengths are superposed.

Coherence. When two beams of light can interact so as to produce interference fringes the visibility of which is unity, they are said to be perfectly coherent. When their interaction produces no fringes (but only photometric summation) they are said to be incoherent or incoherent. An elaborate mathematical theory of coherence recognizes that coherence and incoherence are extreme cases—between them lies “partial coherence.” Zernike, who contributed a great deal to the development of the subject, defined the degree of coherence γ_{12} of two sources as equal to the visibility of the fringes obtained in the most favourable circumstances using light from these sources. It has been shown that the visibility of the fringes obtained in Young’s experiment depends on the width of the slit S_1 , and the following mathematical relation has been derived:

$$\gamma_{12} = \frac{\sin(2\pi ad/l_1)}{2\pi ad/l_1} \quad (10)$$

in which d is the width of the slit S_1 (Figure 4). If d is gradually increased from zero, γ_{12} falls from one (for d equal to zero) to zero for $d = l_1 a/\lambda$, a value equal to the separation of the fringes when $l_1 = l_2$. When the width d is further increased, fringes are again seen but they are of low visibility and are reversed (*i.e.*, there is now a dark fringe in the centre).

For the case of a slit source being inaccessible for measurement, its angular width (d/l_1) can be determined by measuring the visibility of the fringes while a , the separation between P_1 and P_2 , is varied. Michelson used this method to obtain the angular diameter d of a star (serving as the slit source) from measurements of the visibility of interference fringes formed in the focal plane of a telescope that receives light from two small mirrors mounted in front of the telescope’s objective, separated from each other by a distance a .

The concept of coherence that has been applied to light from two pinholes can be extended to a beam of light considered as a whole. A roughly parallel beam of light is incident on a thin sheet of metal normal to the direction of propagation. Then two pinholes may be made in the sheet at A and B (Figure 9) and the visibility of the resulting fringes measured so as to obtain the mutual coherence γ_{AB} . If A and B are initially coincident and are slowly separated then γ_{AB} falls gradually from one to zero. It is possible to define a region of coherence around any point A such that if point B lies within this region the coherence is good ($\gamma_{AB} > 0.7$). Similarly, by devices such as that described in the next section it is possible to measure the mutual coherence between A and a point A' that is, as it were, downstream from A and to define a “coherence length” l_c such that coherence is good when $AA' < l_c$. When, (1) the region of coherence extends across the whole beam of light, and (2) the coherence length is large, the beam is said to be highly coherent because the mutual coherence between any two points such as B and A' is high. What qualifies as a “large coherence length” depends on the type of source and the conditions of the experiment; ten centimetres is a large coherence length for the kind of source considered in the next section, but a well-stabilized gas laser may give a beam with a coherence length of many metres.

In the wave equation for light already cited, the displacement and phase angle, represented by the variables ξ and σ , were used to specify a wave motion, but, for light, these

Partial coherence

Coherence length

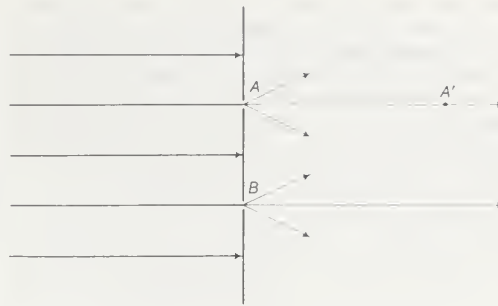


Figure 9: Two pinholes in an opaque sheet to illustrate mutual coherence between points A , A' , and B .

quantities are not observable nor can they be inferred from any observations—because of the high frequency of the wave motion. The coherence γ_{12} and the phase difference, however, are observable quantities that characterize sources and beams of light. This makes them important both in theory and in practice.

Two-beam interference. In the Michelson interferometer, shown in Figure 10, the incident wave W is divided at the beam splitter BS so that part of the light is transmitted and part is reflected. After reflection at M_1 and M_2 the two parts form the wave fronts W_1 and W_2 . These are copies of W , and, because corresponding points are superposed, coherence is obtained even with an extended source. The light from source S , selected by the filter at FF' , is quasi-monochromatic. The plane R represents the image of M_2 that would be seen by reflection in BS . The phase differences between W_1 and W_2 are the same as if W_2 had been reflected from R , which is called the reference plane. M_1 may be traversed normal to itself and may also be tilted with respect to the reference plane. A compensating glass plate C having the same thickness as BS is used so that both wave fronts will pass through a total of three thicknesses of glass.

Fringes will be formed when M_1 is adjusted to be exactly parallel to R and separated from R by a small distance e . For a hollow cone of rays, each ray will be incident on M_1 and on R at an angle θ . After passing through the instrument on their return trip, these rays will be focussed into a circular ring in focal plane FF' of the lens L . At each point on this ring two waves will be superposed and their path difference will be $2e \cos \theta$. Bright rings are obtained for values of θ such that $2e \cos \theta = p\lambda$, in which p is an

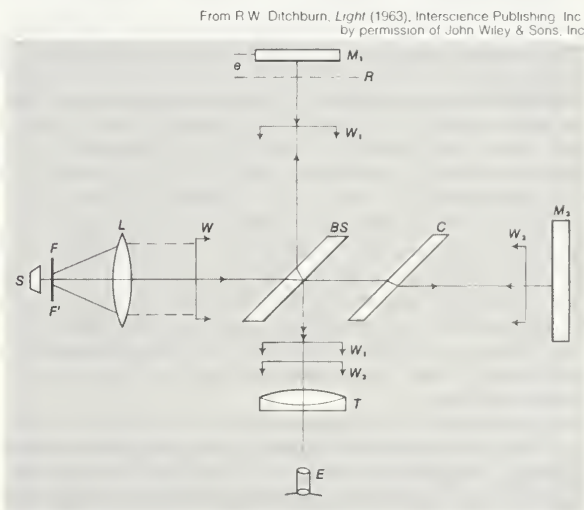


Figure 10: The Michelson interferometer.

integer. The appearance of these fringes is similar to that of Newton’s rings (see Plate).

These fringes are known as fringes of equal inclination because any one ring corresponds to a set of rays that all have the same inclination, θ , to the mirror M_1 . They are conveniently observed by focussing an eyepiece E on the plane F . Because the lens T and the eyepiece E constitute a telescope focussed for an infinite distance, the fringes are said to be localized at infinity.

The apparatus may also be adjusted so that the mirror M_1 is inclined to R , the image plane of mirror M_2 , and nearly in coincidence with it. The incident light is rendered nearly parallel and normal to plane R . If the telescope is removed, straight-line fringes can be seen by an observer who focusses on the region between M_1 and R . A bright fringe is the locus of points for which $2l_p = p\lambda$. These fringes are called fringes of equal thickness.

Fringes of equal thickness may be formed by reflection at the two glass surfaces bounding an air film between two glass plates (Figure 11). Strictly speaking, this arrangement does not give two-beam interference because multiple reflections occur, as shown in the figure. Only the two beams A and B , however, need be considered for present purposes. Beams like B' and C' , caused by multiple reflections, are weak, and, unless the glass plates are fairly thin and of high optical quality, fringes formed by beams reflected from the outer surfaces of the glass plates are close together and are of poor visibility. If the arrangement is such that one of the plates is truly planar and that

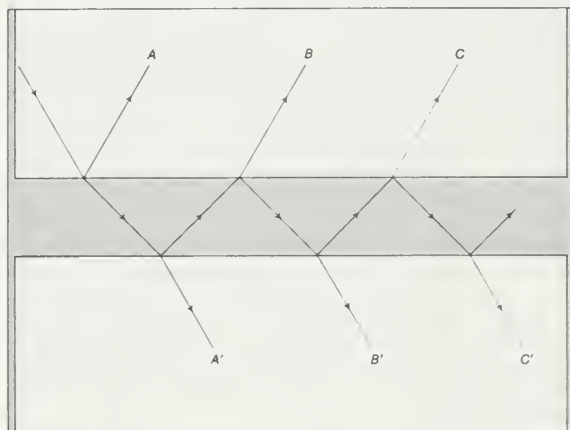


Figure 11: Interference in a thin film of air between two pieces of glass. The rays A, B, C , etc. interfere with each other, as do rays A', B', C' , etc.

the other is spherical, as is the case for a convex lens lying on a glass plate, the resulting fringes of equal thickness are circles centred at the point of contact. They are known as Newton's rings (see Plate). In this situation, in which one surface is plane and the other is not, the fringes form a contour map of the nonplanar surface. They are then called contour fringes. This is a useful method for testing the flatness of a surface and determining the location of irregularities.

Multiple-beam interference. If the two inner surfaces of the plates shown in Figure 11 are coated so as to make them reflect 80 percent or more of the incident light, then the resulting interference pattern will be caused by the superposition of many beams. Figure 12 shows an arrangement for producing the fringes of constant inclination by multiple-beam interference. The amplitudes of successive beams are proportional to r, r^2, r^3 , etc. (r is the ratio of the intensity of the reflected light to that of the incident light for one reflection). The phase differences are $\epsilon, 2\epsilon, 3\epsilon$, etc., in which $\epsilon = (4\pi e \cos \theta)/\lambda$. These fringes are much sharper than those obtained with two-beam interference (see Plate).

With a large number of beams the intensity is extremely high when they are all in phase ($\epsilon = 0$), but, even when the phase difference between any two successive beams (e.g., the first and the second) is quite small, the phase difference between the first and, say, the thirtieth beam is so large that the later beams in the series are in opposition to the earlier beams. Thus the intensity is relatively small except when the value of ϵ is close to one of the values $2r\pi$ (in which p is an integer). Multiple-beam fringes of constant inclination were used by Charles Fabry and Alfred Pérot in France for resolution of spectral lines having only small differences of wavelength. Multiple-beam fringes of constant thickness have been used by an English physicist, Samuel Tolansky, to detect surface irregularities down to less than a nanometre.

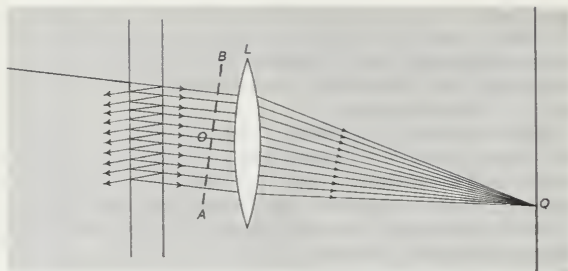


Figure 12: Multiple-beam interference. Lens L concentrates all beams at focus Q with the same phase differences they had while crossing a plane AB normal to OO' .

From R.W. Ditchburn, *Light* (1963), Interscience Publishing, Inc. by permission of John Wiley & Sons, Inc.

Wave groups. If two pendulums that have frequencies ν_1 per minute and $(\nu_1 + 1)$ per minute are started together, they will gradually go out of step; after half a minute they will be moving in opposite directions and after a minute they will be together again. Over a long time they will move together once every minute. In a similar way, when two waves of slightly different frequency are moving in the same direction, they are sometimes in phase and sometimes out of phase so that the resultant is sometimes large and sometimes small, as shown in Figure 13. Two waves may be considered for which the spatial frequencies are ν_s and $(\nu_s + \Delta\nu_s)$ and temporal frequencies ν_t and $(\nu_t + \Delta\nu_t)$. The fluctuation represented by the envelope (dotted line in Figure 13) is called the beat wave. It has a temporal frequency ($\Delta\nu_t$) equal to the difference of the temporal frequencies of the constituent waves and a spatial frequency ($\Delta\nu_s$) equal to the difference of the spatial frequencies. It is therefore propagated with a velocity $U' = \Delta\nu_t/\Delta\nu_s$. Many physical problems involve groups of waves that include a range of frequencies. It is found that, even in a dispersive medium, a group is propagated over a considerable distance as a recognizable unit. The velocity of this recognizable group is $U = d\nu_t/d\nu_s$.

Velocity of groups of waves

From R.W. Ditchburn *Light* 1963, Interscience Publishing, Inc. by permission of John Wiley & Sons, Inc.

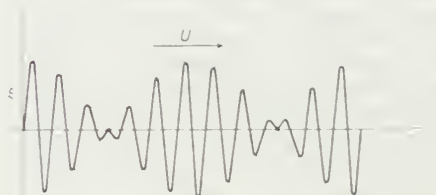


Figure 13: A simple beat wave of amplitude ξ moving with velocity U (see text).

There is a certain kind of wave group for which the variation of the displacement (with distance z) along the path of propagation may be represented by the expression $h(z) \cos \nu_0 z$, in which $h(z)$ is a function that varies with z much more slowly than $\cos \nu_0 z$; e.g., in Figure 13, $h(z)$ would be the function represented by the dotted line, and ν_0 is the spatial frequency of the individual waves represented by the full line. These waves are called modulated waves. If $h(z)$ varies extremely slowly with z , the modulated wave is quasi-monochromatic in the sense described above; i.e., it departs little (over distances long enough to contain many wavelengths) from a monochromatic wave. A modulated wave is completely described when $h(z)$ and ν_0 are known. It is also completely described when the amplitudes and phases of the various waves that make up the group are known. These are given by a function α dependent on the frequency ν , $\alpha(\nu)$. Because $h(z)$ and $\alpha(\nu)$ both describe the same wave group, there must be a relation between them. A mathematical theorem of a French mathematician, Jean-Baptiste-Joseph Fourier, gives this relation, making it possible to calculate either $h(z)$ or $\alpha(\nu)$ when the other is known. The energy density at z is equal to $H(z)$, which is proportional to $[h(z)]^2$. The energy per unit frequency range near ν , is $G(\nu)$, which is proportional to $[\alpha(\nu)]^2$.

Newton's rings

When $h(z)$ varies very slowly with z , $\alpha(v_s)$, and thus $G(v_s)$, is large for a range of v_s close to v_0 and falls rapidly to near zero outside this range, as shown in Figure 14B. If this range in which $G(v_s)$ is large is v_R and if l_R represents a range of z over which $h(z)$ varies very little, then it is found that v_R and l_R are inversely proportional to one another and that their product is of the order of magnitude of unity. This represents the fact that the longer the wave train the more closely its properties agree with those of the ideal monochromatic wave, which is infinitely long and has a precisely defined frequency.

Undisturbed atoms emit exponentially damped waves the length of which is usually 10^7 waves or more so that v_R is

Encyclopædia Britannica, Inc.

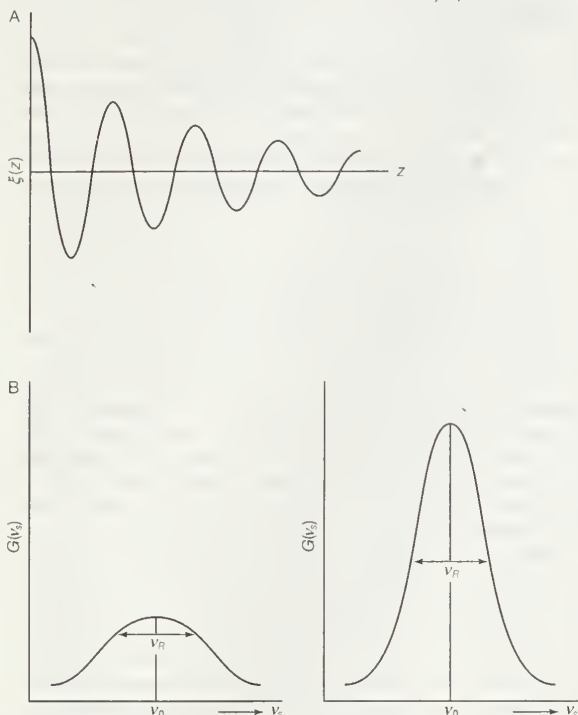


Figure 14: Damped waves. (A) Amplitude, $\xi(z)$, as a function of distance, z . (B) Energy, $G(v_s)$, as a function of frequency, v_s . The figure shows heavy (left) and light (right) damping. For light emitted by free atoms, l_R would encompass 10^7 waves or more and v_R would be correspondingly smaller (see text).

a small fraction of v_0 . Collisions increase the damping by a factor that is proportional to the pressure. The observed radiation is also modified by the Doppler-Fizeau effect, because the atoms that emit the light do not all have the same velocity. This increases the range Δv_R . Even when the effects of collision damping and Doppler-Fizeau effect are combined, the value of l_R for the wave trains emitted in low-pressure electrical discharges is still about 10^5 wavelengths, and most of the energy is confined within a frequency range of order $10^{-5} v_0$ (which corresponds to a wavelength range of less than 0.01 nanometre). These quasi-monochromatic waves are called wave groups. The light that is emitted by high-pressure lamps or by luminescent solids extends over a much wider range of frequency, and wave-group theory has little useful application to problems concerning non-monochromatic light from these sources.

Wave groups in a dispersive medium. In vacuum, all components of the group have the same phase velocity, and therefore the phase relations between different members of the group are constant. A group advances as a unit without any change of the modulation function $h(z)$. In a dispersive medium, the phase relations change and $h(z)$ changes as the wave train advances, but this change is slower than might be expected. Over considerable distances the group is propagated as a recognizable whole with the group velocity U . The change of $h(z)$ is small for passage through a gas and also for groups (for which v_R is small) that represent sharp spectral lines. Thus these wave groups are propagated virtually unchanged through

an optical instrument or, if they arrive from the Sun or stars, through the Earth's atmosphere.

In Young's experiment, Figure 4, the fringes have maximum visibility at the position O , corresponding to zero path difference and zero phase difference for all wavelengths. If a thin sheet of mica is inserted in front of slit P_1 , the centre (or position of maximum visibility) is displaced upward. It was at one time thought that the new centre would be found at the point corresponding to zero phase difference for the mean wavelength in the wave group—*i.e.*, the point calculated for equal times from slits P_1 and P_2 , allowing for the fact that the phase velocity in mica is less than that in air. This did not agree with experimental observation. It was found that the new centre is situated at the position where the times from P_1 and P_2 are equal when the group velocity is used to calculate the time required to traverse the piece of mica. At this position the wave train from P_1 exactly overlaps the wave train from P_2 . At any other position, part of each wave train cannot take part in the interference because it does not coincide with any part of the other wave train. This light that cannot interfere forms a uniform background and so reduces the visibility of the interference fringes.

If white light is used and a fairly thick piece of mica is inserted, no fringes are obtained. This is because the wave train has changed shape so much in passing through the mica that it can no longer match the wave train that has travelled through air.

Michelson, using the apparatus shown in Figure 10, studied the decrease in visibility of interference fringes as the path difference between the two wave trains is increased. The reduction in visibility as the path difference increases may be assigned either (1) to the fact that the parts of the wave trains that overlap are decreasing or (2) to the increasing difference between the positions of the bright fringes for different wavelengths in the group. As was seen above, the length of the wave trains and the range of the wavelength are inevitably linked, and so these alternatives (1) and (2) do not constitute two different theories. They are just two different ways of visualizing how wave trains (or wave groups) interfere.

DIFFRACTION

Theory of diffraction. Huygens assumed that every point on a wave front may be regarded as a source of spherical wavelets the envelope of which is the position of the wave front at a later time. Huygens was thus able to account for rectilinear propagation and for the laws of reflection and refraction. Fresnel added the hypothesis that the wavelets can interfere, and this led to a theory of diffraction. Figure 15 shows how a coherent, monochromatic wave from a point source P falls on the screen

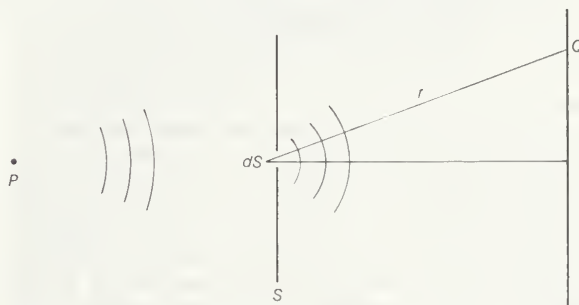


Figure 15: Principle of Fresnel's theory of diffraction (see text).

S , which is opaque except for an aperture dS . Fresnel assumed that the amplitude ($d\xi_Q$) of the wavelet at Q , originating from a small area dS , is:

$$d\xi_Q = \alpha \frac{A}{r} f(\chi) dS, \tag{11}$$

in which A is the amplitude of the incident wave, α is a constant, r is the distance from dS to Q , and $f(\chi)$ is a function of χ , the inclination factor; this factor was introduced by Fresnel because he believed that the effect of the element dS would be greater in the forward direction

Propagation of wave groups from stars

($\chi = 0$) than in an inclined direction. The total effect at Q was obtained by superposing the wavelets from all parts of the aperture, allowing for phase differences caused by a variation of r and also for variation of the inclination factor, $f(\chi)$. Fresnel developed an ingenious method of dividing S into a series of zones of equal area and calculating the total effect as the sum of a simple series. This method applies only to circular apertures and obstacles and then only to points on the axis of symmetry, but Fresnel also developed integrals that are more generally applicable.

Fresnel predicted that there should be a bright spot at the centre of the shadow of a circular obstacle. The experimental verification of this unexpected result gave confidence in Fresnel's wave theory of diffraction.

Fraunhofer diffraction. When the source and pattern screen are sufficiently far from the slit, the phase differences corresponding to different parts dS of the slit opening vary linearly with x and y coordinates in the plane of the aperture (Figure 16). This situation is obtained when two spherical lenses L_1 and L_2 are introduced with source P at the focus of L_1 and Q in the focal plane of L_2 . Spherical waves emanating at the focus of a lens are rendered plane wherever they encounter the lens. Plane waves are made spherical by a lens. They have the same radius of curvature as the focal length of the lens. The wave falling on S is a plane wave, and the total effect at Q may be regarded as caused by a plane wave leaving S . The same result is obtained if L_1 and L_2 are replaced by a single lens (situated near S) that forms an image of P at Q_0 . This

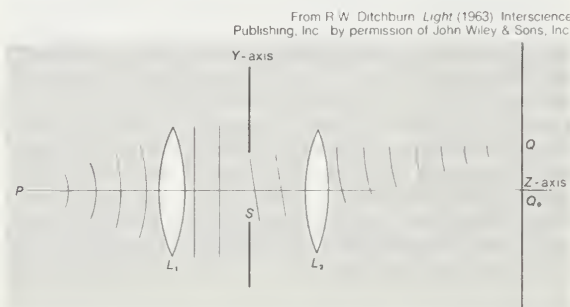


Figure 16: Arrangement for Fraunhofer (far-field) diffraction. The opening at S diffracts light from source P onto plane Q (see text).

is known as far-field diffraction, or Fraunhofer diffraction, and is thus distinguished from near-field, or Fresnel, diffraction. It should be understood, however, that there is only one physical theory of diffraction that is derived from the ideas of Huygens and Fresnel. Fraunhofer diffraction is of great practical importance especially in regard to the performance of optical instruments.

Groups of waves with different directions. When a plane wave is incident upon a slit as shown in Figure 16 (so that its width is limited), the emergent light may be represented by a group of plane waves. All the waves of this group have the same spatial frequency but differ in regard to direction of propagation. It is possible to define a range of angles (in the plane of the page) within which most of the light is found. If this range is θ_R and the width of the slit is w , then it is found that $\sin \theta_R$ is inversely proportional to w ; *i.e.*, the narrower the slit, the greater is the angular spread— $\sin \theta_R$ is roughly equal to λ/d . Using Fourier's theorem it is possible to derive an equation that gives the amplitude and phase of the light diffracted in any direction as a function of the width of the slit. Extension of the calculation to diffraction by apertures or obstacles of any shape involves more lengthy mathematics but no new physical principle.

Angular power spectrum. The energy diffracted in any direction is proportional to the square of the corresponding amplitude. This energy expressed as a function of the angles that define the direction is called the angular power spectrum. It may be measured and is found to agree with that calculated when the width of the slit (or, more generally, the shape and size of the apertures) is known. There are many problems in which it is desired to carry out an inverse calculation—*i.e.*, to calculate the shape and

size of the apertures from measurements of the angular power spectrum. Unfortunately, this is not, in general, possible because measurement of the angular power spectrum does not give the phase of the diffracted light. It is found, however, that measurement of the angular power spectrum yields a function (known as the auto-correlation function) of the size and shape of the obstacles or apertures responsible for the diffraction. In X-ray analysis, this function gives important information about symmetry. A complete picture of the crystal may often be obtained by combining calculations from the angular power spectrum with information derived from other sources.

It is found that when diffraction is due to a number of apertures (or obstacles) that are similar in size, shape, and orientation, the angular power spectrum (G) is the product of two factors, F and f , in which F (called the form factor) depends only on the properties of the individual aperture and f (called the structure factor) depends only on the arrangement or spacing of the elements. When the apertures are irregularly arranged, f is just equal to N (the number of apertures). Thus the diffraction halos produced by an irregular distribution of small similar objects have the same intensity distribution as the pattern for a single particle. This principle is used in a device called an eriometer to determine the size of blood corpuscles and may also be used to calculate the average size of the small particles that cause a halo around the Moon.

When N similar elements are arranged in a regular pattern, the structure factor may vary from zero to N^2 . A diffraction grating (a plate having parallel lines engraved across its surface) with N lines is such a pattern, and, for any of the directions θ_p defined by equation (6), the light from all elements (lines) is in phase; the amplitude is N times that given by a single element, and thus the energy and the structure factors are proportional to the total number of lines squared—*i.e.*, $f = N^2$.

Limits of resolution. Diffraction spreads the light in optical images; so that if two objects are too close to each other, the gap between them cannot be distinguished. The distribution of intensity with radius in the image of a point source is shown in Figure 17. Rayleigh showed, theoretically and experimentally, that the images of two point sources are just resolved when their separation is such that the centre of the pattern due to one image falls on the first minimum of the pattern due to the other (Figure 18). This implies that a telescope with a perfect objective of diameter D can just resolve two stars whose angular separation is $1.2 \lambda/D$. Qualitatively, this agrees with a calculation that shows that most of the energy in the diffraction pattern of an aperture of width d lies within an angular range $\pm \lambda/d$.

The angular separation of the maxima resulting from light

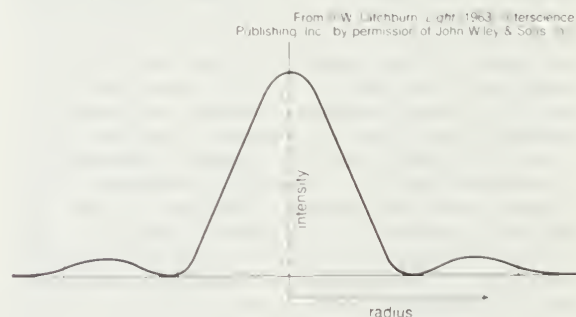


Figure 17: Illumination of a point source image modified by diffraction, shown as the variation of intensity with radius

of two wavelengths λ and $\lambda + \Delta\lambda$ in a spectrum formed by a diffraction grating is obtained by differentiating equation (6), resulting in $\Delta\theta = p \Delta\lambda / c \cos \theta$. These maxima are just resolved if $\Delta\theta = \lambda/M$, in which M is the width of the beam diffracted by the grating. For a grating of N lines, this width is $Nc \cos \theta$, and the resolving power $R = \lambda / \Delta\lambda = pN$. A grating ten inches wide, for example, with 10^4 lines per inch and $p = 10$, has a resolving power $R = 10^6$.

The limit of resolution for a microscope depends on conditions of illumination and is at best about half a wavelength ($\lambda/2$), or about 250 nanometres for visible light.

Diffraction by similar apertures

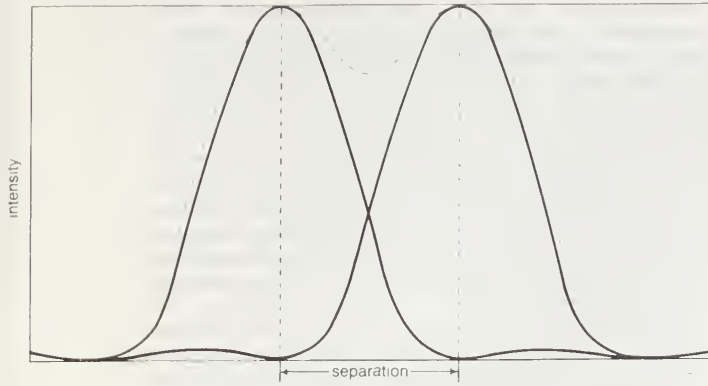


Figure 18: Overlapping images of two point sources. Full lines show how intensity varies with distance from a separate source, dashed line shows combined intensity.

From R.W. Ditchburn, *Light* (1963), Interscience Publishing, Inc., by permission of John Wiley & Sons, Inc.

Polarization and electromagnetic theory

POLARIZED LIGHT

Interaction of plane-polarized beams. Fresnel and Arago, using an apparatus based on Young's experiment (Figure 4), investigated the conditions under which two beams of plane polarized light may produce interference fringes. They found that: (1) two beams polarized in mutually perpendicular planes never yield fringes; (2) two beams polarized in the same plane interfere and produce fringes, under the same conditions as two similar beams of unpolarized light, provided that they are derived from the same beam of unpolarized light or from the same component of a beam of unpolarized light; (3) two beams of polarized light, derived from perpendicular components of the same beam of unpolarized light and afterwards rotated into the same plane (*e.g.*, by using some device such as an optically active plate) do not interfere under any conditions.

Result (1) is to be expected because two displacements in perpendicular planes cannot annul one another, and result (2) is also easily understood. Result (3) shows that mutually perpendicular components of unpolarized light in a beam are non-coherent. Their phase difference varies in time in an irregular way. Unpolarized light has a randomness, or lack of order, as compared with polarized light (implying an entropy difference). This order (or lack of order), rather than the azimuthal property, is the most fundamental difference between polarized and unpolarized light. Perfectly monochromatic light is perfectly coherent and completely polarized.

Superposition of polarized beams. Two coherent beams of plane polarized light may be thought of as propagated in the Oz direction, one with its vector along Ox and the other with the electric vector along Oy ; *i.e.*, the two vibrations are at right angles to each other as well as to the direction of propagation (Figure 19). If the beams have amplitudes a_x and a_y and phases ϵ_x and ϵ_y , then, in general, the resultant vibration (R_x , R_y , and R_z) may be represented in magnitude and polarization by a vector, or arrow, the tail of which touches the axis of propagation Oz while the point moves round the ellipse (Figure 19). It

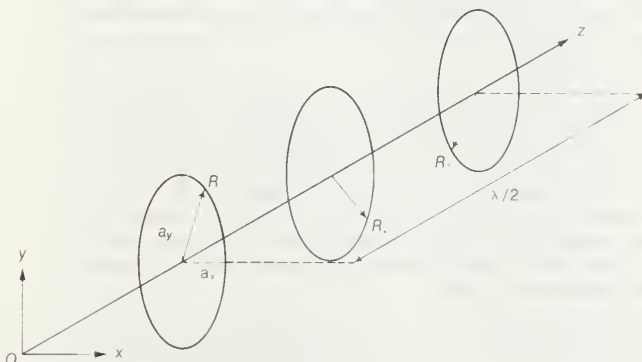


Figure 19: Progression of elliptically polarized wave (see text).

Interference of polarized light

goes round once when the phase angle ϕ (see equation [1]) changes by 2π —*i.e.*, at any given place when t changes by τ or for any one time when z changes by λ . The beam is said to be elliptically polarized. If the phase difference is $\pi/2$, then the axes of the ellipse are equal to a_x and a_y and are along Ox and Oy .

Elliptically polarized light may be regarded as the most general type of polarized light. If the amplitudes of the two waves are equal, $a_x = a_y$, and the phase difference is still $\pi/2$, the ellipse becomes a circle and the light is said to be circularly polarized. If the phase difference ϵ_{xy} is not equal to $\pi/2$, the resultant is still elliptically polarized light, but the axes of the ellipse no longer coincide with the axes of coordinates. If the phase difference $\epsilon_{xy} = 0$ or π , the ellipse shrinks to a straight line and the light is said to be plane-polarized. If the representative vector, when viewed by an observer who receives the light, rotates in a clockwise direction, the light is said to be right-handed (or positive) elliptically polarized light. The opposite sense of rotation corresponds to left-handed (or negative) elliptically polarized light.

In the above analysis, elliptically polarized light is regarded as the resultant of two beams plane-polarized in perpendicular planes. Conversely, it is possible to regard plane-polarized light as the resultant of two beams of

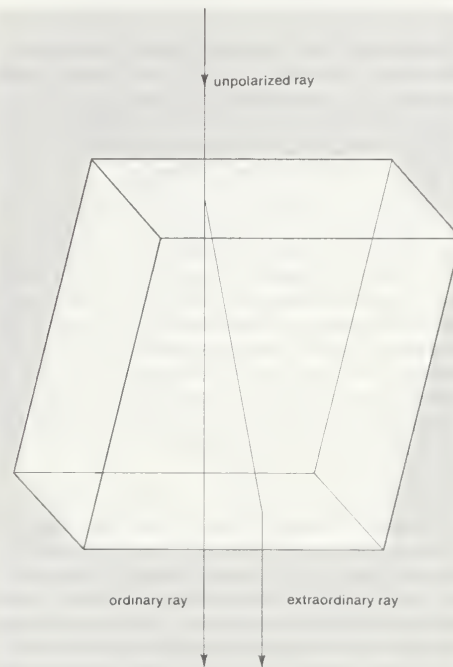


Figure 20: Double refraction showing two rays emerging when a single light ray strikes a calcite crystal at right angles to one face (see text).

elliptically (or circularly) polarized light of the same wavelength, provided that the ellipses are similar in orientation and eccentricity, but one beam is right-handed and the other left-handed.

Double refraction. In the 17th century Bartholin showed that a ray of unpolarized light incident on a plate of calcite, unlike glass or water, is split into two rays, as shown in Figure 20. One ray, called the ordinary ray, is in the plane containing the incident ray and the normal to the surface. If the angle of incidence is varied, this ray is found to obey Snell's law of sines, equation (3). The other ray, called the extraordinary ray, is not in general coplanar with the incident ray and the normal; also, for it, the ratio of sines is not constant. The fact that Snell's law is not obeyed in certain directions implies that the velocity of light in such a medium, called anisotropic, depends on the direction of travel in it. The two rays are polarized in mutually perpendicular planes. This is known as double refraction, or birefringence.

In order to apply Huygens' method of constructing wave fronts (see above *Theory of diffraction*), it is necessary to assume that, in an anisotropic medium, the wave surface

Birefringence

from a point source consists of two sheets, or surfaces (Figure 21). The observation that one ray obeys both laws of refraction implies that one sheet must be a sphere—like the wave surface in an isotropic medium. Huygens assumed that the other sheet is an ellipsoid of revolution that touches the sphere either internally (Figure 21A) or externally (Figure 21B). There is one velocity of propagation for the direction defined by the line through the two points of contact (called the optic axis) and two velocities for any other direction corresponding to light polarized in two mutually perpendicular planes.

From R.W. Ditchburn *Light* (1963) Interscience Publishing, Inc. by permission of John Wiley & Sons, Inc.

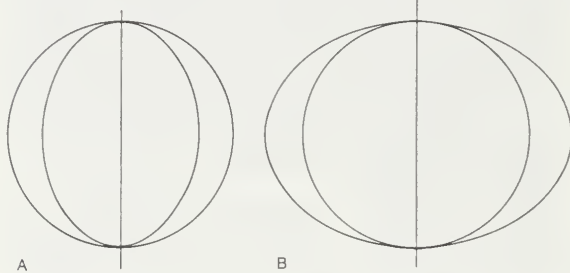


Figure 21: Wave surface for (A) positive and (B) negative uniaxial crystals.

Crystals for which the wave surface has the form shown in Figure 21A are said to be positive uniaxial crystals, and Figure 21B refers in a similar way to negative uniaxial crystals. Huygens thought all crystals were uniaxial, but later observations showed that the general form of the wave surface is more elaborate and is biaxial.

When a parallel beam of plane-polarized light is incident normally upon a thin crystal plate and the crystal is rotated about an axis normal to the plate, two orientations will be found in which a single beam of plane-polarized light emerges; for one orientation, the light is called the ordinary ray, and for the other orientation, the extraordinary ray. Two lines may be drawn on the plate (or on its mount) to indicate the direction of the electric vector of the incident beam when the orientation is such that plane-polarized light emerges. The directions of these lines are called privileged directions for the given anisotropic plate; they are perpendicular to one another.

When light waves pass through an isotropic plate of thickness d , the phase change is $2\pi nd/\lambda$ (n is the index of refraction). For an anisotropic plate there are two indices, n_o and n_e , corresponding to the two privileged orientations, and there are corresponding phase changes, $2\pi n_o d/\lambda$ and $2\pi n_e d/\lambda$. If a beam of light, polarized in some plane other than one of the two privileged orientations, is incident on the plate, it may be resolved into two wave components (lying in two mutually perpendicular planes) that emerge with a phase difference ϵ_{oc} equal to $2\pi(n_e - n_o)d/\lambda$. The phase difference ϵ_{oc} is called the retardation of the plate. When the retardation ϵ_{oc} is equal to $\pm\pi/2$, the plate is called a quarter-wave plate, and when ϵ_{oc} is equal to $\pm\pi$, it is called a half-wave plate.

The maximum difference of ($n_e - n_o$) for calcite, a crystal that is strongly birefringent, is about 0.17. Quarter-wave plates are often made by splitting sheets of mica. Inasmuch as the relevant difference of indices is about 0.004, a mica quarter-wave plate is about 60 wavelengths, or 0.03 millimetre thick.

When a parallel beam of plane-polarized light is incident normally on a quarter-wave plate, the emergent light is circularly polarized if the plane of the electric vector bisects the angle between the privileged directions. Otherwise the emergent light is elliptically polarized, the axes of the ellipse being parallel to the privileged directions of the quarter-wave plate. In a similar way, a suitably oriented quarter-wave plate may be used to convert elliptically polarized light to plane-polarized light.

Many substances that are normally isotropic become weakly birefringent when subjected to shear stress (photoelastic effect; see Plate). Birefringence is also induced by an applied electric field (Kerr effect). A restricted class of crystals that are normally birefringent shows a large change

Induced
birefringence
effects

in birefringence when an electric field is applied. Weak birefringence is produced by an applied magnetic field.

Production of polarized light. Unpolarized light may be separated into two components, polarized in perpendicular planes, either (1) by reflection or (2) by refraction. A single reflection at the polarizing angle produces a well-polarized reflected beam of low intensity (*i.e.*, most of the waves are polarized in one plane) and a strong transmitted beam that is only partially polarized. A pile of glass plates provides many surfaces for reflection and thus gives a much stronger reflected beam that is still well polarized; it also considerably improves the degree of polarization of a transmitted beam.

A Nicol prism (Figure 22) is made by cutting a calcite crystal in a suitable plane relative to the crystal axes and cementing the two parts together with Canada balsam. The critical angle for total reflection is less for the ordinary ray than for the extraordinary ray. The crystal is cut at an angle (relative to the crystal axes) such that, for a cone of incident rays of divergence about 24° , the ordinary ray is totally reflected, whereas the extraordinary ray is transmitted. Both rays are almost completely polarized. A prism of this type with air instead of Canada balsam as the separation film was devised by Foucault for polarization of ultraviolet radiation.

From R.W. Ditchburn *Light* (1963) Interscience Publishing, Inc. by permission of John Wiley & Sons, Inc.

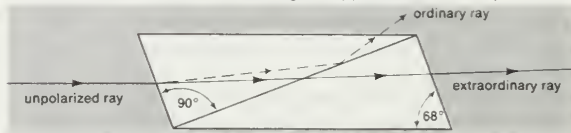


Figure 22: Nicol prism (see text).

Many crystalline substances have different absorption coefficients for the ordinary and extraordinary rays, but the difference is not large enough to provide a useful method of polarization. A series of artificial materials that polarize by absorption has been developed. Complex crystals are produced in a plastic matrix and oriented by stretching. One material of this type, called Polaroid H, transmits about 80 percent of plane-polarized light with the electric vector in one plane and less than 1 percent if the vector is in a perpendicular plane. This type of film offers the most convenient method of obtaining an intense beam of light that is nearly 99 percent plane-polarized.

Analysis of polarized light. Any polarizer has a preferred plane such that if an incident plane-polarized beam has the electric vector in this plane, a high intensity is passed, whereas the device rejects light polarized with the electric vector in a plane perpendicular to the preferred plane. This property implies that any polarizer may be used as an analyzer to test whether or not a given beam of light is plane-polarized. The analyzer is rotated in the beam and if the intensity I of the transmitted light is zero for some orientation, the beam is plane-polarized with the electric vector in a plane perpendicular to that for which the intensity is zero. If there is a position for which the analyzer gives a non-zero minimum, then either (1) the light is elliptically polarized or (2) it is partially polarized. Also if the intensity is not affected by rotating the analyzer, then the light is either (3) unpolarized or (4) circularly polarized. These ambiguities between (1) or (2), and (3) or (4), may be removed and the analysis of a mixture may be completed by using an analyzer and a quarter-wave plate together.

Stokes parameters and Poincaré sphere. Any mixture of different kinds of polarized light may be reduced to an equivalent amount of one kind of polarized light characterized by the orientation and length of axes of an ellipse. It may, of course, happen that the ellipse degenerates to a straight line or a circle. An alternative way of characterizing polarized light is by means of the Stokes parameters, named after a British physicist, George Gabriel Stokes, and defined as follows:

$$\begin{aligned} S_0 &= \langle a_x^2 \rangle + \langle a_y^2 \rangle; & S_1 &= \langle a_x^2 \rangle - \langle a_y^2 \rangle; \\ S_2 &= 2\langle a_x a_y \cos \epsilon_{xy} \rangle; & S_3 &= 2\langle a_x a_y \sin \epsilon_{xy} \rangle, \end{aligned}$$

in which the polarized light is regarded as the resultant

of two beams polarized with vectors a_x and a_y along the coordinates Ox and Oy , and ϵ_{xy} is the phase difference as in the section above. The brackets $\langle \rangle$ indicate time averages. If the light is completely polarized, then the square of S_0 is equal to the sum of the squares of the other three parameters. This relation was used by a French mathematician, Henri Poincaré, to produce an elegant representation of the properties of polarized light by means of a sphere, which is known as the Poincaré sphere. Each point on the sphere represents a different kind of polarized light; e.g., the two poles represent right- and left-handed circularly polarized light, points on the equator represent plane polarized light; other points on the sphere represent different kinds of elliptically polarized light. The Stokes parameters constitute a matrix representation of polarized light. This kind of representation makes available certain powerful mathematical methods (which have been developed in other connections) for calculating how the state of a beam of polarized light changes on reflection or on passing through a crystal.

Colours of thin plates. When a parallel beam of white light (I_0) is passed through a thin slice of crystal placed between a polarizer and an analyzer, the intensity I of the emergent light may be expressed as the sum of two terms:

$$I = \text{white term} + \text{colour term.}$$

The white term represents the amount of light that would be transmitted if the crystal were not present. It is equal to $I_0 \cos^2 \theta$, in which θ is the angle between the preferred directions of polarizer and analyzer. The colour term depends on the relation between the directions of the crystal axes and the preferred directions of polarizer and analyzer. It is dependent on the retardation ($\epsilon_{oc} = 2\pi d [n_e - n_o]/\lambda$) of the plate as well. Because this retardation varies with the wavelength, the transmission depends on the colour. Some parts of the spectrum are transmitted more than others, and the emergent light is coloured. The first term, on the contrary, does not depend on the properties of the crystal, and it represents white light. It can be reduced to zero by setting the analyzer and polarizer so that their preferred directions are at right angles ($\theta = \pi/2$). The colour is then obtained most strongly, because any transmission of light results from anisotropy (i.e., birefringence) in the slice. This provides a sensitive method of detecting strains in glass and may also be used to detect induced anisotropy when mechanical or electrical stress is applied. This method is sometimes used to detect regions of large strain in beams or girders. A model of the beam is made in plastic and placed between polarizer and analyzer. On applying a load, fringes appear in the regions of most strain. Plate 1 shows a picture obtained in this way. The example is chosen for simplicity, whereas useful application to more complicated structures requires laborious calculations.

A high-speed optical shutter called the Kerr shutter can be made by placing an isotropic material between crossed polaroids, one acting as a polarizer and the other as an analyzer, and applying an electric field. If a suitable material is used, the intensity of the transmitted light will reproduce the variation of electric field, even at frequencies over a megahertz. An improved shutter is produced by using the change of birefringence in some crystals.

Rings and brushes. The retardation of waves by a crystal plate depends on the relation between the perpendicular to the surface of the plate and the axes of the crystal from which it was cut. Observation of the colours of one plate with parallel light gives information in relation to one direction within the crystal. It is possible to obtain information about many directions within the crystal at one and the same time. This is done by taking the light that has passed through the polarizer and converging it strongly, using a microscope objective (reversed). It then passes through the crystal as a solid cone of light. It is rendered nearly parallel again by a second objective and passes through the analyzer. Finally it passes through an optical system that brings all rays that have passed through the crystal in one direction to a focus at a single point in a certain plane (a similar focussing is shown in Figure 12). Thus each point in this plane corresponds to

a certain direction in the crystal and to a single value of retardation ϵ_{oc} . The differences in retardation ϵ_{oc} are caused partly by variation of $(n_o - n_e)$ with the direction in the crystal and, usually to a lesser extent, by a difference in the length of path within the crystal. This implies that, with monochromatic light, bright and dark fringes are observed. Each dark fringe is the locus of points for which the retardation is an odd number multiple of 180° —i.e., ϵ_{oc} equals $p_0\pi$ (p_0 being an odd integer). If the slice is cut from a uniaxial crystal with its normal in the direction of the optic axis, the symmetry of the wave surface (Figure 21) requires that the fringes be circular rings (see Plate). The rings are interrupted at certain points for which the second term (in the expression of the intensity of the light transmitted by a crystal plate placed between polarizer and analyzer) is zero. This gives either a bright cross or a dark cross depending on the angle between the preferred directions of the polarizer and the analyzer; these crosses are known as brushes. The patterns obtained with a slice cut so that the normal is not in the direction of the optic axis (a crystal direction for $n_o = n_e$) or with biaxial crystals are complicated and sometimes beautiful (see Plate). Such patterns may be used to identify crystals present as fairly small inclusions in other materials.

ELECTROMAGNETIC-WAVE CHARACTER OF LIGHT

Maxwell's equations. Historically the theory of electricity and magnetism developed in the form of a number of empirical laws each of which was a generalization based on a series of experiments; e.g., Coulomb's law dealt with the force between two stationary electric charges. Maxwell replaced all these laws by a single theory concisely stated in the form of a set of vector equations. It has been said that Maxwell's theory is Maxwell's equations, and indeed it is impossible to do justice to Maxwell's achievement without use of these equations. The content of these equations and their relevance to the theory of light will be described here in general terms. The reader may refer to standard texts on electromagnetic theory for the equations themselves.

Electric and magnetic fields are specified by means of the vectors E and H with which are associated the vectors D , B , and J (electric and magnetic induction and density of electric current). Maxwell's equations fall into two groups: (1) three constitutive equations and (2) four field equations.

All material bodies contain electrons. These are negative charges circulating around heavier nuclei that are positively charged. When an electric field is applied to a material body, the average positions of the negative charges relative to the positions of the positive charges are changed. This creates an internal electric field. Similarly the action of a magnetic field on a material changes the movement of the electrons and sets up an internal magnetic field. The constitutive equations state that effects within a material body are proportional to the applied fields so that the resultant fields within the body are proportional to the applied field. Certain constants are defined: ϵ is the dielectric constant, μ is the magnetic permeability, and σ is the electrical conductivity. There is no general agreement, however, concerning these constants, and therefore some authorities use a different nomenclature (for treatment of the constants in greater detail, see ELECTROMAGNETIC RADIATION AND ELECTRICITY AND MAGNETISM). But these are basic properties of the material: for free space the values are ϵ_0 , μ_0 , and zero respectively.

The first of the four field equations quantifies certain properties of the electric induction at the boundary of a volume that contains a net positive or negative charge. The second states that, since there are no free magnetic poles, a certain integral of magnetic quantities is zero. The third equation states that, when the magnetic flux through a surface changes, electrical voltages appear on the boundaries of the surface. The fourth states that electrical currents in conducting materials and changes of the electric induction in nonconducting materials produce magnetic effects.

The field equations (like the constitutive equations) are linear equations: they state that certain quantities are proportional to one another; e.g., the third equation states

Constitutive and field equations

The Kerr shutter

that the electrical voltages are directly proportional to the rate of change of magnetic flux. No new constants (other than ϵ , μ , and σ), however, are introduced. In Maxwell's equations, electricity and magnetism are two aspects of one thing called electromagnetism. Maxwell's theory can indeed be stated in a way that does not mention electricity and magnetism separately.

Maxwell's
new
hypothesis

The three constitutive equations and the first three of the field equations are precise formulations of known empirical laws. In the fourth field equation Maxwell introduced a new hypothesis—that an electrical change in a nonconductor produces magnetic effects. This hypothesis—which can be verified by electrical experiments—leads to the theory of electromagnetic waves capable of being propagated through a vacuum.

Propagations of electromagnetic waves. When Maxwell's equations are combined, using standard mathematical methods, a new equation is obtained. This is similar in form to the wave equation (7) given above. It predicts the existence of electromagnetic waves, with well-defined properties which will now be described.

1. In free space (*vacuo*) electromagnetic waves are propagated with a phase velocity $c = (\epsilon_0\mu_0)^{-1/2}$; the magnitude of c obtained from electrical measurements of $\epsilon_0\mu_0$ is approximately equal to the measured velocity of light (see Table 1). Plane waves are propagated without attenuation. Spherical waves have an amplitude inversely proportional to the distance from a small source.

2. For transparent nonconducting mediums such as water or glass, the phase velocity is $b = (\epsilon\mu)^{-1/2}$, and the index of refraction is $n = c/b = (\epsilon_0\mu_0/\epsilon\mu)^{-1/2}$. This relation holds for radio waves and, with moderate accuracy, for infrared radiation. For visible light it holds moderately well for some mediums but not for others. The constants ϵ and μ are associated with redistribution of charge in response to changes in the electromagnetic field. At low frequencies these movements follow the changes of the field in a simple way. At optical frequencies, when the field is reversing about 10^{14} times a second, the situation is more complicated. The effective value of μ at these frequencies is μ_0 and the effective value of ϵ is less than the value measured with static fields, or fields that vary more slowly—apart from the possibility of resonance, which will be considered later.

3. In conducting mediums, the electromagnetic waves are absorbed as they are propagated. This absorption is exponential; *i.e.*, the fraction absorbed in a thin layer of thickness d is βd , in which β is a constant the value of which can be calculated when ϵ and σ are known. For metals, β is very large so that a layer only a tenth of a micrometre in thickness absorbs about half of the incident light. This absorption is much stronger than that of nonconducting substances, which are normally regarded as opaque (such as ebony); these are usually quite transparent in thicknesses of a micrometre or so.

Although the simple theory is correct in predicting high absorption for metals in general, the values of β for different metals (calculated from ϵ and σ) are not correct. This is because the effective value of σ at high frequencies is not the same as that measured with direct current or low frequency alternating current.

4. The waves are transverse; the electric field vector (E) and the magnetic field vector (H) are perpendicular to one another and to the direction of propagation along the axis OZ as shown in Figure 23.

5. The two quantities E and H , the magnitudes of vectors E and H , fluctuate in phase, so that when E is a maximum H also is a maximum (Figure 23). This is implied in the equation

$$\epsilon^{1/2}E = \mu^{1/2}H, \quad (12)$$

which enables either E or H to be calculated when the other is known.

6. The electromagnetic field possesses energy, and the energy per unit volume (W) is equal to $1/2(\epsilon E^2 + \mu H^2)$. In view of equation (12) above, this implies that for electromagnetic waves,

$$W = \epsilon E^2. \quad (13)$$

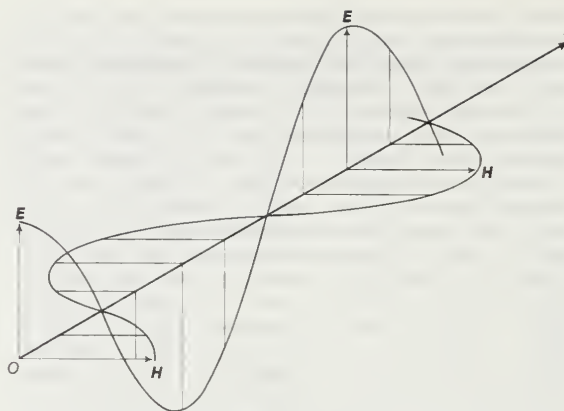


Figure 23: Electromagnetic wave, showing that electric field vector E and magnetic field vector H are in phase (see text).

7. The propagation of electromagnetic waves involves a transport of energy. A vector $S = E \times H$, called the Poynting vector (after an English physicist, John Henry Poynting), is used to describe this flow—*i.e.*, the energy flows in the direction of S (which is also the direction of wave propagation in an isotropic medium)—and the amount of energy crossing unit area (lying across the direction of flow) in unit time is equal to S —*i.e.*, to the magnitude of S . For electromagnetic waves $S = E \times H = (\epsilon/\mu)^{1/2}E^2$, and this relation, combined with equation (13), implies that $S = W/(\epsilon\mu)^{1/2}E^2$, or $S = bW$; thus, the mean value of S is b times the mean value of W —*i.e.*, the energy crossing unit area in unit time is equal to the velocity multiplied by the energy density.

The
Poynting
vector

The linearity of Maxwell's equations implies that the law of superposition applies to electromagnetic waves. The electromagnetic theory is therefore able to give a satisfactory account of all the simpler phenomena of interference and diffraction already described above. There are a few problems for which it is clearly superior to the "scalar-wave" theory that was described. These include the polarization of light transmitted through a slit the width of which is of the same order of magnitude as the wavelength and also the polarization of light scattered from microscopic particles. Calculations based on the electromagnetic theory are sometimes difficult mathematically, but the results are in agreement with experiment.

The interaction of light with matter

REFLECTION AND REFRACTION

When light is incident in a perpendicular direction on the surface of a piece of glass, part is transmitted and part is reflected. This happens at any surface that forms the boundary between two transparent mediums of different refractive indices. For a particle theory of light, this phenomenon constitutes a serious difficulty. If a beam of light consists of a stream of particles all alike, because like causes produce like effects, there is no good reason why some should go through the surface and some be turned back. If the particles are not all the same (*e.g.*, if some have more energy and are therefore able to get through), then it would be expected that all those particles that do penetrate one surface would be able to pass through a second. This is never observed. In fact, the fraction of light transmitted at the second surface is the same as the fraction transmitted at the first (for light of one wavelength). Newton was aware of this fundamental difficulty and introduced the ad hoc hypothesis of particles with periodic oscillation between "fits of easy reflection" and "fits of easy transmission."

A wave theory of light, on the other hand, offers an explanation that follows in a natural and logical way from its premises. When a wave passes from one medium to another certain conditions must be fulfilled at the boundary. These conditions arise because the frequency of the transmitted wave must be the same as that of the incident wave. The velocities of light in the two mediums are different, and therefore, because velocity is the product of

wavelength times frequency, the wavelength cannot be the same. By way of analogy, elastic waves can be imagined falling on a welded boundary between two solids—e.g., steel and brass—with point P_1 just inside the first medium and P_2 extremely near P_1 and just inside the second medium. Then if P_1 and P_2 do not move exactly together the material at the boundary must tear. If the boundary is so strong that it does not tear, then the waves at the boundary must match each other in a most precise way. This match is possible only if there is a reflected wave as well as a transmitted wave: when the density and the elastic properties of the two mediums are known then the fraction reflected can be calculated.

For light, the fraction reflected can be calculated when the indices of refraction are known, and it is found that the fraction reflected when a beam of light is incident perpendicular to a boundary between air (refractive index nearly equal to 1) and a medium of refractive index n is $(n-1)^2/(n+1)^2$. Thus for glass (index 1.5) about 4 percent is reflected, but for diamond (index 2.42) 17 percent is reflected. This result is predicted by any reasonably self-consistent wave theory including the electromagnetic-wave theory. When light is incident on a surface at an angle other than perpendicularly, however, the calculation of the fraction of light reflected is more difficult. This fraction depends both on the direction of incidence and on the polarization of the light.

Fresnel relations

Fresnel derived certain expressions (known as Fresnel's relations) for the amount of light reflected under various conditions, using an elastic solid theory. To do this he assumed boundary conditions chosen ad hoc—i.e., to produce the correct result. In the electromagnetic theory there is no choice in regard to the conditions that must be satisfied at the boundary of two dielectrics. The boundary conditions are fixed by the constitutive equations and by general principles such as the conservation of energy. They are derived from experiments on electricity, not from experiments on light.

These boundary conditions may be applied to electromagnetic waves with the following results:

1. The existence of a wave reflected so that the angle of reflection is equal to the angle of incidence and of a refracted wave the direction of which obeys Snell's law.

2. All the results that Fresnel had previously obtained for the amounts of light reflected and refracted under a variety of conditions, including Brewster's law.

3. The existence of total reflection and of certain observed results in regard to the change of polarization of light that has been totally reflected (when the incident light is polarized).

4. The amount of light reflected at the surfaces of metals (including the result that strong reflection is, in general, obtained with good conductors).

5. A detailed description of certain complicated results obtained when light is incident upon the surface of a transparent anisotropic crystal (e.g., calcite).

DISPERSION AND SCATTERING

Oscillating dipoles. In electrostatics, a dipole consists of two equal and opposite charges situated a small distance apart. The line joining the charges is called the axis of the dipole, the dipole moment M is the name given to the product of one charge times separation. An oscillating dipole may be thought of as one in which charges move so that M fluctuates between $+M_0$ and $-M_0$ (i.e., $M = M_0 \sin 2\pi\nu t$). Maxwell's theory shows that an oscillating dipole produces electromagnetic waves that are nearly spherical at distances large compared with the size of the dipole. The total energy emitted is proportional to the square of the ratio M_0/λ^2 (in which $\lambda = c/\nu$).

Scattering by free electrons. When an electromagnetic wave (with an electric vector of amplitude E_0) operates on a free electron, it causes the electron to oscillate so as to produce an oscillating dipole with a moment proportional to E_0/λ^2 . This dipole emits waves in all directions. These waves are regarded as scattered light. The scattered light is strongest in directions perpendicular to the axis of the dipole, which is the same as the direction of the electric vector in the incident wave. It is polarized with its electric

vector in the plane containing the dipole. The ratio of the energy scattered to the energy incident per unit area is: $k = (1/6\pi\epsilon_0^2) (e^4/m^2c^4)$, in which e is the charge and m is the mass of the electron. The energy scattered is equal to the amount of energy incident upon an area equal to k . For this reason k is called the scattering cross section or the scattering coefficient.

Scattering by bound electrons. A bound electron is subject to a restoring force when displaced from its equilibrium position, similar to a weight on the end of a spiral spring. An electron bound to an atom can oscillate with a natural frequency ν_0 . Its motion produces an oscillating dipole and electromagnetic waves are emitted. Energy lost in this way constitutes damping, and the oscillation decays like the well-known damped harmonic oscillator of classical mechanics. The curve of decay is similar to the curve of a damped wave shown in Figure 14A. When a light wave is incident, a bound electron oscillates. The strength of the dipole produced and, hence, of the light scattered depends on the frequency (ν_i) of the light wave and also upon the frequency ν_0 and the damping constant γ of the electron.

Dipole radiation

The scattering cross section (k') is found to be:

$$k' = \frac{\nu^4}{(\nu_0^2 - \nu_i^2)^2 + \gamma^2\nu_0^2} k; \quad (14)$$

and it reaches a maximum value when ν_i is nearly equal to ν_0 ; i.e., when the incident wave is in resonance with the natural frequency.

Molecular scattering. Because atoms and molecules contain bound electrons, they scatter light, but not just like single bound electrons. The scattered wave has an amplitude proportional to the amplitude (E_0) of the incident wave and inversely to the distance from the scattering centre; i.e., the electric vector E_r of the scattered wave is $\eta E_0/r$, in which η is a constant. This value applies to the wave emitted at right angles to the electric vector of the incident light. The phase of the scattered wave is either the same as that of the incident wave or differs from it by an amount that is the same for each molecule. Figure 24 shows a beam of light incident upon a layer of gas at low pressure. The

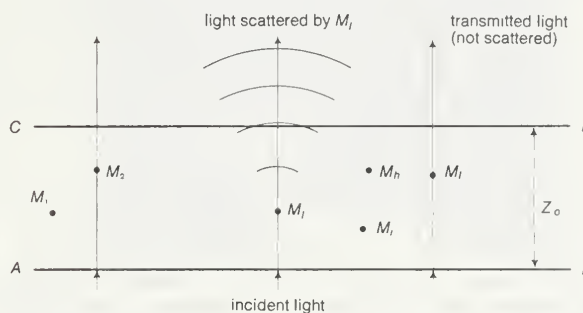


Figure 24: Scattering of light from molecules (see text).

molecules are irregularly placed so that the light scattered from different molecules, such as M_1, M_2, M_n, M_i , etc., in any direction (except the forward direction) is incoherent. The law of photometric summation (see above) applies, so that the total energy scattered by N molecules is just N times the energy scattered by one molecule.

A gas with N molecules per unit volume has an equivalent cross section (k_m),

$$k_m = \frac{8\pi}{3} \eta^2 N. \quad (15)$$

In the forward direction the light scattered from different molecules is coherent because the total distance from the back edge of the gas layer (AB) to any molecule, such as M_j and M_i and thence to the front edge of the layer (CD) is equal to the thickness of the layer (z_0). Thus in the forward direction, the scattered waves can interfere with each other and with the wave representing the unscattered light. The resultant is a new wave in the forward direction but retarded—i.e., with a phase lag $\eta\lambda Nz_0$. In ordinary

wave theory, however, the velocity of light in a medium of refractive index n is c/n and the phase lag due to passage through a layer of thickness z_0 is equal to $2\pi(n-1)z_0$. Equating these two estimates of the lag requires

$$2\pi(n-1) = \eta\lambda^2 N, \quad (16)$$

and substituting for η from equation (16) into equation (15) implies that

$$k_m = \frac{32\pi^2}{3\lambda^4 N} (n-1)^2. \quad (17)$$

This equation was derived by Rayleigh. He also showed that the brightness of the blue sky is about that to be expected if the light is produced only by molecular scattering, and it is not necessary to suppose that there is any appreciable scattering by dust particles—in the high atmosphere and on a clear day. The factor $1/\lambda^4$ implies that the scattering coefficient for blue light ($\lambda = 400$ nanometres) is about six times larger than that for red light ($\lambda = 640$ nanometres) and this accounts for the blue colour. Laboratory measurements of the amount of scattered light agree with that calculated from equation (17). For the rare gases (*e.g.*, argon) the scattered light is almost completely plane-polarized; but other gases show less complete polarization (*e.g.*, 90 percent for carbon dioxide) as would be expected if their molecules possess an electric dipole moment. The light of the blue sky is polarized but polarization is reduced because some of the light has been scattered more than once.

Dispersion. Equation (17) implies that $(n-1)$ may be calculated if η is known; *i.e.*, if the ratio of the amplitude of the scattered wave to that of the incident wave can be calculated. Dispersion formulas for gases have been obtained by assuming that the atoms or molecules contain bound electrons having a response that varies with frequency in the manner described above. The methods of calculation involve the properties of atoms and molecules rather than those of light and are not considered in this article. Certain general results are readily obtained—*e.g.*, that for most transparent mediums the refractive index increases regularly toward the blue end of the spectrum (this is called normal dispersion). It is also found that in a region in which there is a strong absorption, the index varies rapidly, and over part of this region the dispersion is anomalous (*i.e.*, the index increases in the direction of larger wavelengths). This derives from equations (14) and (17).

Scattering of larger particles. A particle that contains N atoms but is small compared with the wavelength of light has a dipole moment approximately N times that for a single atom. The scattered energy is then proportional to N^2 , or to V^2 , or to R^6 (if V is the volume, and R is the radius). Thus a particle 20 nanometres in diameter scatters about as much light as 10^{12} separate atoms. As the size of the particle increases, the scattering continues to increase, though more slowly. When the particle is larger than a wavelength it must be considered as a diffraction object. The ratio of effective cross section for scattering to actual cross section fluctuates as the size increases. The theory of this variation, in relation to size, refractive index, and absorption has been worked out and verified experimentally.

Scattering by macroscopic solids and liquids. In crystalline solids the atoms are regularly arranged so that the wavelets scattered in a given direction have regularly varying phases. They interfere to give a small resultant except in the direction of the forward wave. Scattering of light is small in relation to the number of atoms per unit volume and results largely from crystal defects. For X-ray wavelengths smaller than the interatomic distance, the coherent waves scattered in certain directions reinforce one another.

For amorphous solids and liquids the number of atoms in any volume having a diameter of about one wavelength is large and the scattering per atom, although usually larger than in crystals, is still small in relation to the number of atoms involved.

The Kramers–Kronig relation. Any satisfactory theory of dispersion must comply with the condition that the scattered wave can never appear in advance of the incident wave that produces it. Physicists Hendrick Anthony Kramers of The Netherlands and Ralph de Laer Kronig of Germany showed that this basic causality condition implies that the dispersion (*i.e.*, the variation of refractive index with frequency) and the absorption are not independent. They derived equations enabling the absorption to be calculated when the dispersion is known (for all frequencies) and vice versa. It is not surprising that a relationship should exist, because dispersion and absorption are each related to the resonators described above in connection with scattering by bound electrons. The relationship has been found of great importance in many branches of pure and applied physics.

MECHANICAL EFFECTS OF LIGHT

When light is emitted or absorbed there are three mechanical effects: (1) an exchange of energy, (2) an exchange of linear momentum, and (3) an exchange of angular momentum. These are manifested as (1) heating (or cooling), (2) a pressure, and (3) a torque. If the conservation laws apply, then light must possess energy, momentum, and angular momentum. All these effects are predictable either from classical electromagnetic theory or from quantum theory. They can all be explored in relation to either exchanges between light and macroscopic pieces of matter or interaction between radiations and atoms. In this section the first will be treated, leaving the latter for the next section.

Equation (13) is an expression for light energy according to electromagnetic theory, and every detector of light is an energy transducer (converter). There is a general theorem, valid for light and sound, predicting that radiation must exert a pressure equal to W (the energy per unit volume) on a body that absorbs it. Twice this pressure is produced if the light is totally reflected at normal incidence. In electromagnetic theory, when light falls at normal incidence on a metal surface, a current (proportional to the field E) is produced in the metal; and, because this occurs in the magnetic field of the light, there is a reaction proportional to the vector product of the electric and magnetic field strengths ($E \times H$) and normal to the surface—*i.e.*, a reaction proportional to the Poynting vector (S). For an insulator the induced current is the displacement current. Because the Poynting vector is proportional to W , the energy per unit volume, the pressure is proportional to W (further consideration shows that, in the usual units, it is equal to W).

The measurement of the linear momentum of light is difficult because, under normal experimental conditions, the pressure is less than 10^{-4} dyne per square centimetre. In Figure 25, light from powerful sources S_1 and S_2 falls

From R.W. Ditchburn *Light*, 1963, Interscience Publishing, Inc. by permission of John Wiley & Sons, Inc.

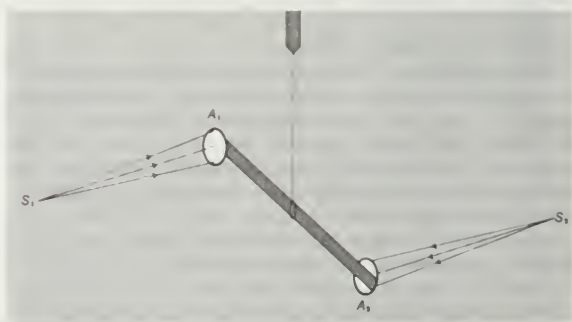
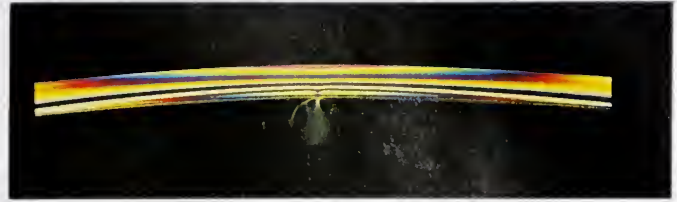
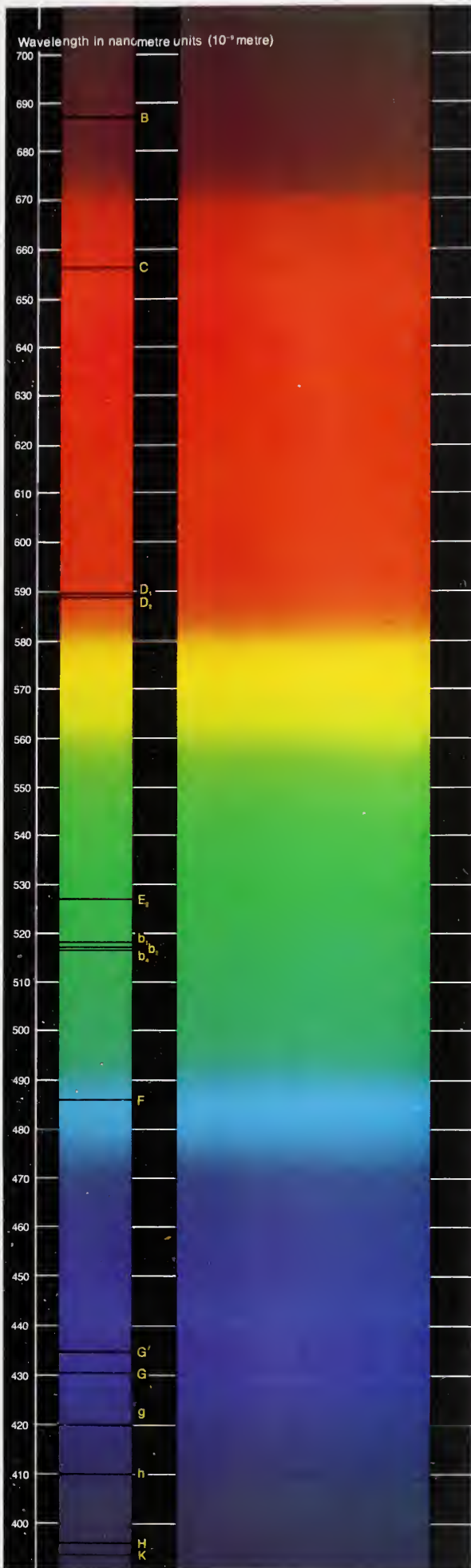


Figure 25: Apparatus for measurement of light pressure (see text).

on the mirror vanes A_1 and A_2 that are at the ends of a rod suspended by a quartz fibre. The pressure on the vanes produces a couple or twisting of the fibre and hence a rotation of the system that is measured by the usual mirror-and-scale method. The true radiation pressure can be measured only when a certain thermal action, called the radiometer effect, has been eliminated. In the radiometer

Normal and anomalous dispersion

Light pressure

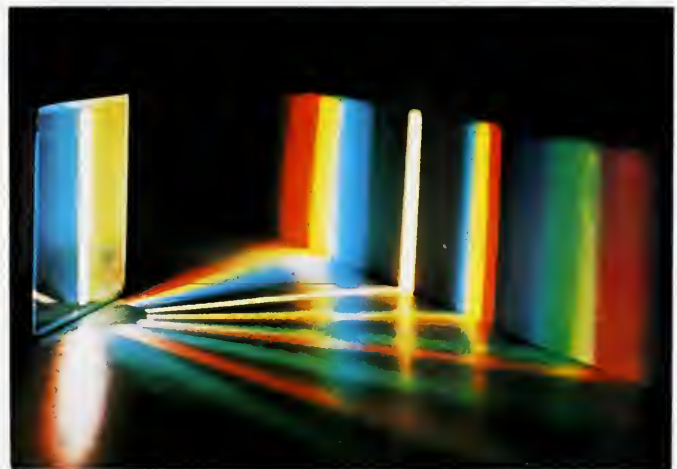


Bending of transparent plastic bar between crossed Polaroids to show photoelastic strain patterns.

The visible solar spectrum (simulated), showing prominent Fraunhofer absorption lines.



Spectrum of white light by (above) a prism and (below) a diffraction grating. With a prism, the red end of the spectrum is more compressed than the violet end.



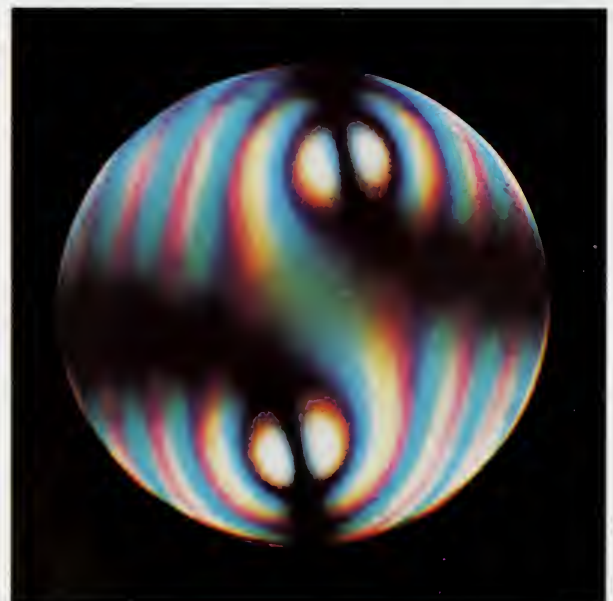
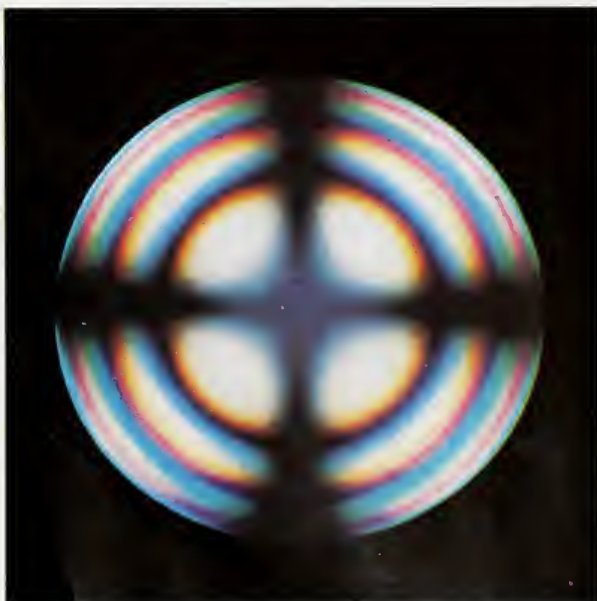


Fabry-Pérot fringes.

Newton's rings.



Interference patterns obtained with crystals in convergent plane-polarized light (polarizer and analyzer crossed): (left) uniaxial crystal cut normal to the axis; (right) biaxial crystal cut perpendicular to bisector of the axes.



effect, if radiation heats a surface, then molecules of gas that rebound from one side will move faster than those that rebound from the other, cooler side, and there will be a net pressure. In the period 1901 to 1905, when physicists were making the first accurate determination of light pressure, it was not possible to remove this effect by placing the system in a vacuum in which the number of molecules is at a minimum; instead, the effect was reduced by substituting the special cells for the vanes. Any radiometer effect within the cells produces no net force. The residual effects at outer surfaces are small and are the same whichever side of the cell is illuminated. The measured values (1) when radiation is absorbed and (2) when it is reflected agreed with theoretical values within 1 percent. By use of an alternative method of eliminating the radiometer effect, results in agreement with the theoretical values were also obtained.

It can be shown that a beam of circularly polarized light possesses angular momentum and thus should exert a torque upon a half-wave plate that reverses the rotation of polarization (*i.e.*, from left-handed to right-handed circularly polarized light or vice-versa). The experimental difficulties of measuring this torque (of about 10^{-11} dyne centimetre) are formidable and the overall error is about 10 percent. The result agrees with the theoretical value within this margin.

The density of radiation incident on the Earth's surface at its mean distance from the Sun is about 1.36 kilowatts per square metre (for a surface normal to the Sun's rays), and the corresponding light pressure is about 5×10^{-5} dyne per square centimetre. The total pressure on the Earth is about 6×10^{14} dynes, a minute fraction of the gravitational force exerted by the Sun.

Importance of the mechanical properties of light

The mechanical properties of light, difficult to measure in the laboratory, are important under certain astronomical conditions. Gravitational forces are proportional to the mass of a body and hence to the cube of its radius. For a body having large radius compared with the wavelength, the radiation pressure is proportional to the square of the radius. Whereas the attraction and repulsion may be equal and balanced for a medium size particle, for a smaller particle, the radiation pressure may be a large fraction of the whole force. For this reason, radiation pressure has a considerable effect in extending a comet's tail. In a star the outward flow of radiation from the centre produces a net force tending to expand the star, thus opposing the gravitational forces. Also, when radiation energy is transferred between two parts of a rotating star that are moving with different velocities, momentum transfers are involved. This produces an effect known as radiative viscosity because the resulting forces reduce the relative velocity of adjacent layers and so reduce the angular velocity of rotation.

Quantum theory of light

In the classical electromagnetic wave theory, the absorption of light is a continuous process and there is no lower limit to the amount of energy that an atom can absorb from light of given frequency. Energy can be exchanged, in infinitesimal amounts between radiation and atoms, in an enclosure that is in thermal equilibrium. The total energy of the enclosure should, in equilibrium at temperature T , be distributed so that each degree of freedom has an equal amount, kT , in which k is Boltzmann's constant. The number of possible modes of vibration for transverse waves with frequencies between ν and $\nu + d\nu$ is $8\pi\nu^2 d\nu$ per unit volume. (In this section all frequencies are temporal frequencies and ν is used for ν_t). If $\rho(\nu)d\nu$ is the radiation energy per unit volume in this frequency range, and c is the velocity of light, then the radiation density is as shown in curve R-J of Figure 26. This law, derived by Rayleigh and Jeans, is the inevitable consequence of classical wave theory. It predicts that the radiation density increases without limit as the frequency increases, and this implies that, in equilibrium, all the energy in the universe is found in the high-frequency end of the electromagnetic spectrum and none is present in matter. This catastrophe may be avoided by assuming that there is some

upper limit to the possible frequency of radiation, but the Rayleigh-Jeans law does not agree with measurements of temperature radiation in the visible and ultraviolet regions of the spectrum (curve R-J, Figure 26). The experimental results are represented by curve P in Figure 26, and Planck suggested a formula, now called Planck's law, that fitted this curve. Originally this formula was probably an inspired guess, but Planck soon showed that it could be derived by assuming that a mode of frequency can change only by $h\nu$ (or by a multiple of $h\nu$) in which h is a new universal constant now called Planck's constant. The minimum "element of energy" (as Planck referred to it) is now called the quantum. For visible light the quantum is about 3.5×10^{-12} erg or 3.5×10^{-19} joule.

The theory in which Planck introduced the quantum was concerned with the interaction between radiation and matter. He expected that only minor modifications of the electromagnetic theory of light in free space and of the electron theory of matter would be needed. This expectation was not realized. Experiments, some of which will be described below, led Einstein to regard light as an assembly of entities, later called photons. A frequency ν , energy $h\nu$, momentum $h\nu/c = h/\lambda$, and angular momentum $h/2\pi$ are associated with a photon of circularly polarized light. In many situations the photons appear to be localized

Properties of a photon

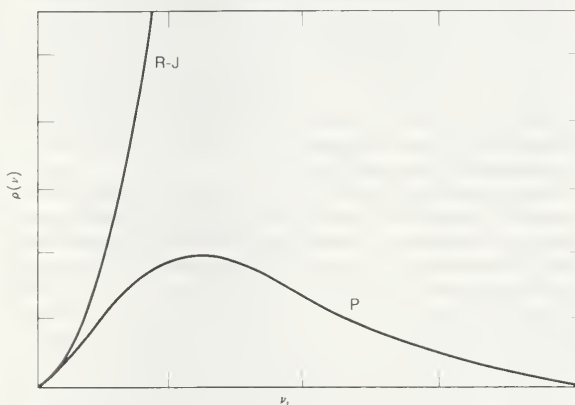


Figure 26: Equilibrium radiation density curve R-J (Rayleigh-Jeans); curve P (Planck; see text).

concentrations of energy and momentum—*i.e.*, to have the properties of particles. Yet the evidence that led to the rejection of a corpuscular theory in favour of a wave theory is still as valid and as compelling as it was in the time of Fresnel.

The theory, which includes logical descriptions both of those experiments that seem to require a corpuscular theory and those that support a wave theory, is known as quantum mechanics. This is a general theory of radiation and matter—not a special hypothesis about their interaction. The account given below of some aspects of quantum mechanics that are important in relation to experiments on light must be related to the general theory (see MECHANICS: *Quantum mechanics*).

PHOTONS

When light of suitable wavelength is incident upon the clean surface of a metal, electrons are emitted. If W_m is the maximum kinetic energy of one of those photoelectrons released by light of frequency ν , then it is found that the maximum kinetic energy is Planck's constant times the difference between two frequencies, or $W_m = h(\nu - \nu_0) = h\nu - W_0$, in which W_0 is written for $h\nu_0$. No electrons are emitted when ν is less than the critical frequency ν_0 . It is known from other experiments that a minimum amount of energy called the work function is needed for an electron to escape from a surface. If W_0 is accepted as the work function for emission of a photoelectron, then it is easily understood that no electrons are emitted when frequency ν is less than the critical frequency and that the maximum energy of escape for an electron is $h\nu - W_0$. Some electrons may make collisions after absorption of energy $h\nu$ and either fail to emerge from the surface or else emerge with lower energy. Escape

by gaining energy from two photons of frequency less than ν_0 is extremely improbable because the energy due to one absorption is almost always lost by collision before a second photon is absorbed.

The number of photons absorbed, and therefore the number of electrons released, are proportional to the energy density in the beam of light and to the time. The energy of individual electrons is independent of the energy density—*i.e.*, independent of the amplitude of the waves.

Conservation of energy and momentum. Niels Bohr, a Danish physicist, postulated that an atom can exist only in one of a set of discrete stationary states of internal energy (W_1, W_2, \dots). The frequency $h\nu_{nm}$ of light emitted in passing from state n to m is related to the energy difference by Bohr's frequency condition: that Planck's constant times the frequency is equal to the energy difference between two states; *i.e.*,

$$h\nu_{nm} = W_n - W_m.$$

For some pairs of states ($W_n - W_m$) can be measured by experiments on collisions between electrons and atoms, and Bohr's condition is found to be satisfied within the limits of experimental error. Many other, less direct but more accurate, experiments are also in agreement with this equation.

In Bohr's theory, and subsequently in quantum mechanics, an angular momentum is associated with each atomic state. In the simplest situation, when polarized light is emitted in a magnetic field, the angular momentum about a known axis changes by $h/2\pi$. The changes are such that angular momentum is conserved if a photon of circularly polarized light has an angular momentum of $h/2\pi$ about the direction of propagation.

It was shown by a United States physicist, Arthur Holly Compton, that when X-radiation is scattered by a nearly free electron, there is a wavelength change. A formula that gives the wavelength change as a function of the angle of scattering was derived by assuming that a "particle" of energy $h\nu$ and momentum $h\nu/c$ collides with a particle of mass m that is at rest, and that energy and linear momentum are conserved. The resting electron gains momentum and kinetic energy by changing its speed. The photon cannot change its speed, which is always c . It therefore must change its frequency when it changes energy and momentum. The measured change of frequency agrees with that calculated by Compton.

An atom that has absorbed a quantum of energy $h\nu_{nm}$ and changed from state m to state n may return to state m with emission of a photon of the same frequency (resonance radiation)—*i.e.*, $\nu_{mn} = \nu_{nm}$. If there is a state s intermediate in energy between states m and n , the return may be in two stages; n to s and s to m , with the emission of two photons. The sum of the two quanta $h\nu_{ns}$ and $h\nu_{sm}$ is then equal to a single quantum $h\nu_{nm}$ so that energy is conserved. In gases, these emissions usually follow the absorption in less than a microsecond, and the term fluorescence is used to include both resonance radiation and multistage processes. In solids and liquids there are sometimes intermediate states that are metastable so that emission occurs over much longer periods (phosphorescence; see below *Luminescence*). In all these processes the frequency of the emitted radiation is not greater than that of the absorbed radiation in accord with an empirical law discovered in the 19th century. In quantum theory this law is a necessary consequence of conservation of energy.

In Rayleigh scattering, the frequency of the scattered light is equal to that of the incident light. The physicists Sir Chandrasekhara Venkata Raman of India and Leonid Isaakovich Mandelstamm of the Soviet Union independently discovered processes in which light is scattered by an atom or molecule that changes its state. There is a corresponding discrete change in the frequency of the scattered radiation. When the scattering atom or molecule is in its lowest state, the frequency of the unscattered radiation is always less than that of the incident light (Stokes lines). When light is scattered by a molecule that is not in its lowest state, photons of higher frequency may be found in the scattered radiation (anti-Stokes lines). The scattered photon then has more energy than the incident photon,

but the molecule gives up energy by passing to a lower state. The Raman effect is observed when light is scattered by solids and liquids. There is also an additional kind of scattering in which photons exchange energy and momentum with the thermal vibrations of the solid or liquid (Brillouin scattering, named after a French-U.S. physicist, Léon Brillouin).

The experiments quoted are examples of a large number that are in accord with the view that, within the limits of measurement, energy and momentum are conserved in detail—*i.e.*, in each individual interaction between an atom and radiation.

Spontaneous and stimulated emission. Planck's theory gives the distribution of radiation energy with frequency in an assembly of atoms and radiation in thermodynamic equilibrium at a common temperature in an enclosure. An Austrian physicist, Ludwig Boltzmann, applying the same principles to the assembly of atoms, derived a formula that gives the distribution of atoms between different stationary states—*i.e.*, the ratio of the number N_n in any other state n to the number N_m in any other state m . In the enclosure now considered, Boltzmann's and Planck's laws must both be satisfied though each atom is changing its state frequently by emitting or absorbing a photon or by collision. Equilibrium can be maintained if, and only if, each individual reaction is balanced by its own reverse process. Einstein assumed that the equilibrium was maintained by the following process: (1) absorption of radiation at a rate $B_{mn}N_m\rho(\nu)$; (2) spontaneous emission at a rate $A_{nm}N_n$; and (3) stimulated emission at a rate $B_{nm}N_n\rho(\nu)$, in which A_{nm} , B_{mn} , and B_{nm} are constants, now known as Einstein coefficients. When the rate of absorption was equated to the rate of emission (*i.e.*, to the sum of spontaneous and stimulated emission), an equation of the same form as Planck's law was obtained. This equation became identical with Planck's law if

$$B_{nm} = B_{mn} = \frac{c^3}{8\pi h\nu^3} A_{nm}. \quad (18)$$

If there were no stimulated emission (*i.e.*, if $B_{nm} = 0$), Planck's law and Boltzmann's law could not both be satisfied. Einstein postulated the existence of stimulated emission in 1905, but for more than 30 years it appeared to be only of theoretical interest. At the densities of radiation then available, the rate of stimulated emission was small compared with the rate of spontaneous emission. Equation (18) shows that in the microwave region of the electromagnetic spectrum (frequencies of a few thousand megahertz) the ratio B_{nm}/A_{nm} is much larger than in the visible spectrum. Accordingly, it was much easier to develop microwave amplification by stimulated emission of radiation (maser) than the corresponding laser effect with radiation in or near the visible spectrum. The existence of stimulated emission can be fitted into a classical wave theory, like many other concepts that originated in relation to theories of photons.

Interference and diffraction of photons. It was at one time suggested that interference and diffraction phenomena are caused by collisions between photons considered as small particles, or at least to some kind of complex interaction between photons. A simple experiment performed by a British physicist, Geoffrey Ingram Taylor, in 1908 excludes this possibility. He photographed the diffraction pattern of a needle and reduced the illumination until long exposures were needed to obtain an image. When the chance of two or more energy quanta passing through the apparatus simultaneously was made extremely small, the diffraction pattern was exactly the same as that obtained with a strong source of light.

This experiment is supported by many later experiments, showing that interference and diffraction are to be associated with single photons. The discussion in the section on coherence above implies that when ordinary sources of light are used each photon can interfere only with itself and not with any other photon. An interference pattern can, of course, be photographed only by recording the effects of many photons because one photon can activate only one grain on the photographic plate, but it does not

Bohr's
postulate

Compton
scattering

Raman
effect

Einstein
coeff-
icients

matter whether the photons all arrive over a time span of a microsecond or of several weeks.

The photograph of an interference pattern is both a wave phenomenon because it shows the characteristic spatial periodicity and a quantum phenomenon because the whole energy of a photon can be used to activate a single grain. It is not possible to trace the path of one photon (regarded as a particle) through the apparatus and, at the same time, obtain the interference pattern. In Young's experiment there is no way of finding out through which slit a given photon passes—except by covering one slit and thus losing the interference fringes.

THE WAVE-PARTICLE NATURE OF LIGHT

Uncertainty relations. Experiments with the photoelectric effect show that energy can be transferred from one atom to another in a way that suggests that photons are corpuscles—*i.e.*, localized concentrations of energy and momentum. Other experiments imply equally clearly that the light emitted from an atom, when it loses energy, must be represented by wave groups. Both of these sets of experiments are equally valid, and together they require that the photon must have both wave-group and particle properties at the same time.

It will now be assumed that the intensity of the waves gives the relative probability that a suitable detector will detect a photon at a given point. Experimentally, waves and particles are both abstractions, each describing the same physical system. Quantum mechanics does not seek a relationship such that particles are guided by waves according to the laws of classical physics—as in the theory of the English physicist J.J. Thomson, in which photons are pictured as closed tubes of electric force guided by Maxwellian waves. The waves and particles now to be considered do not have all the relevant properties of classical waves and particles. When a photon is absorbed, the whole of a wave group, possibly extending over a large volume, is annihilated. The position of a photon particle cannot be exactly specified. It is known only that it is within the wave group and most likely exists where the waves are most intense. Also the energy and momentum of the photon are determined by the frequency of the waves. Fourier analysis of wave groups shows that, when the wave group is short (so that the position of the photon is fairly well determined), the frequency of the waves and, hence, the energy and momentum of the photon are known only within a wide range. The product of the uncertainty in position and the uncertainty in momentum can never be less than a certain minimum value. It will now be shown how these relations appear in some experimental situations.

In one experiment a parallel beam of light is incident on a slit of width d , producing a diffraction pattern on a screen. The y -coordinate, at the moment when a photon passes through the slit, is known within an uncertainty $\Delta y = \pm d$. The direction of propagation when the photons emerge from the slit is uncertain because of diffraction, but most of the photons are observed within an angular range $\pm\theta$, in which $\sin \theta = \lambda/d$. There is therefore an uncertainty Δp_y in the y -component of the momentum (p) of about $\pm(h \sin \theta)/\lambda$, so that the product of the two uncertainties, Δp and Δy , is

$$\Delta p \Delta y \sim \frac{dh \sin \theta}{\lambda} \sim h,$$

in which the sign \sim indicates is of the "order of magnitude of." In a second experiment, monochromatic light is passed through an optical shutter. The shutter is opened for a short time Δt so that a wave group of duration Δt and length $\Delta z = c\Delta t$ passes through the shutter. Then $c\Delta t$ is the length l of the wave train. This implies that the frequency ν varies within a range $\Delta \nu = 1/\Delta t$ —*i.e.*, that the photon energy varies within a range $\Delta E \sim h/\Delta t$ so that $\Delta E \Delta t \sim h$. More detailed analysis shows that, if the uncertainties Δt and ΔE are the root mean square deviations of statistical theory, then their product is equal to or larger than $h/4\pi$; so that $\Delta E \Delta t \geq h/4\pi$.

Individual examples do not prove that the uncertainty

relations are generally true any more than individual examples prove the general validity of the second law of thermodynamics. The fact that every example agrees, and that ingenious attempts to find exceptions all fail, constitutes a kind of proof by default, but the real proof is in the whole range of experimental work that demonstrates the equal status and inescapable association of wave and particle properties.

Maximum observations. The uncertainty relations show that precise measurement of both position and momentum are incompatible. It follows that a precise deterministic theory of particle dynamics is not appropriate because the precise initial condition of a particle is not observable. An extension of the uncertainty relations shows that the number of photons in a light beam and its phase cannot be simultaneously measured with precision, and thus a classical wave theory is also excluded.

In any experimental situation there are, in general, two or more observations that are compatible. Neither affects the other, and their results are independent of the order in which they are made. A complete assembly of such measurements is said to constitute a maximum observation. A state of a system is specified by giving the result of a maximum observation or in an equivalent way. The relative probabilities of future observations on a system can be calculated from the result of a maximum observation without reference to any previous results obtained on the system. For every system there is one certain prediction—that an immediate repetition of the maximum observation will give the same result.

Empirical questions. There is a certain type of question to which the quantum mechanics gives no direct answer because the question is not related to a possible observation. One may ask, for instance, what happens to a photon of right-handed circularly polarized light as it passes through a half-wave plate. It is known that the photon will emerge as a photon of left-handed circularly polarized light and that angular momentum is transferred to the plate. The wave theory gives an account of the progress of a wave through the plate. If, however, one tries to visualize the particle and asks for a mental picture of the way in which this particle alters as it goes through the plate, there is no ready answer. The theory has nothing to say unless an experiment is designed to determine what happens to the photon in passage through the plate. For example, a half-wave plate is split in two and the light is passed through the quarter-wave plate, which forms one half. It may subsequently be passed through a Nicol prism analyzer. There will be some orientation of the Nicol prism such that the photon will certainly pass through. If a Nicol prism is oriented at α to this direction, the probability for a single photon to pass through is $\cos^2 \alpha$. This detailed theoretical prediction agrees with experiment. Thus the quantum mechanics correctly predicts the results of an experiment but leaves unanswered the original question about the way in which the photon, visualized as a particle, changes its angular momentum as it passes through the half-wave plate.

Assemblies of photons. The application of quantum mechanics to the electromagnetic field in an enclosure justifies a representation in terms of photons. These may be regarded as indistinguishable particles. Any number of these particles may exist simultaneously in the same energy state. The statistics of particles that obey these conditions is called Einstein-Bose statistics. The distribution of particles among different states may be calculated by the usual methods of statistical mechanics in which probability (or entropy) is maximized subject to certain limitations. In an assembly of material particles the total energy and also the number of particles are constants. With photons, the second condition is no longer valid because two photons of frequencies ν_n and ν_m may be replaced by one of frequency $\nu_{nm} = \nu_n + \nu_m$, thus decreasing the number of photons while conserving the total energy. The statistical calculation leads straightforwardly to Planck's law. It is found that for ordinary sources of light the probable number of photons in any given state is always much less than unity. For a laser source many photons are found in relatively few states, and, indeed, in an ideal single

mode continuously operated laser, all photons would be in one state.

A considerable amount of experimental and theoretical work has recently been done on coincidences of photon detection in two or more detectors that are observing the same source. This kind of photon statistics measures the coherence of light from a star and hence the size of the star.

Assemblies of photons and atoms. In the Dirac theory of absorption and emission of radiation, quantum mechanics is applied to an enclosure that contains radiation and free atoms or molecules. Processes involving creation and annihilation of photons with corresponding atomic transitions must maintain equilibrium. The theory derives expression for the Einstein *A* and *B* coefficients of absorption and emission and for their ratios. It is also shown that, whereas spontaneous emission is random, stimulated emission is coherent with the radiation that excites it. The theory extends to include an account of Rayleigh scattering (and hence of dispersion) and of the Raman effect. All problems concerned with the absorption, emission, and scattering of light are, in principle, soluble, although the mathematical formulation, even of approximate formulas, is usually difficult.

The quantum-mechanical theory is adequate at optical frequencies and at lower frequencies. It is also adequate at the higher frequencies associated with ultraviolet radiation and for X-ray scattering including the Compton effect. It is inadequate for photons of the still higher frequencies involved in the creation and annihilation of fundamental particles. It is to be expected that, in due course, a basic theory of electromagnetic forces and of their relation to gravitation and to the strong forces between nucleons will emerge. The present theory will then have the status of an approximation valid for a wide range of frequencies including the frequencies associated with visible light.

(R.W.Di.)

Luminescence

SOURCES AND PROCESS

Luminescence is the emission of light by certain materials when they are relatively cool. It is in contrast to light emitted from incandescent bodies, such as burning wood or coal, molten iron, and wire heated by an electric current. Luminescence may be seen in neon and fluorescent lamps; television, radar, and X-ray fluoroscope screens; organic substances such as luminol or the luciferins in fireflies and glowworms; certain pigments used in outdoor advertising; and also natural electrical phenomena such as lightning and the aurora borealis. In all these phenomena, light emission does not result from the material being above room temperature, and so luminescence is often called cold light. The practical value of luminescent materials lies in their capacity to transform invisible forms of energy into visible light.

Luminescence emission occurs after an appropriate material has absorbed energy from a source such as ultraviolet or X-ray radiation, electron beams, chemical reactions, and so on. The energy lifts the atoms of the material into an excited state, and then, because excited states are unstable, the material undergoes another transition, back to its unexcited ground state, and the absorbed energy is liberated in the form of either light or heat or both (all discrete energy states, including the ground state, of an atom are defined as quantum states). The excitation involves only the outermost electrons orbiting around the nuclei of the atoms. Luminescence efficiency depends on the degree of transformation of excitation energy into light, and there are relatively few materials that have sufficient luminescence efficiency to be of practical value.

Luminescence and incandescence. As mentioned above, luminescence is characterized by electrons undergoing transitions from excited quantum states. The excitation of the luminescent electrons is not connected with appreciable agitations of the atoms that the electrons belong to. When hot materials become luminous and radiate light, a process called incandescence, the atoms of the material are in a high state of agitation. Of course, the atoms of

every material are vibrating at room temperature already, but this vibration is just sufficient to produce temperature radiation in the far infrared region of the spectrum. With increasing temperature this radiation shifts into the visible region. On the other hand, at very high temperatures, such as are generated in shock tubes, the collisions of atoms can be so violent that electrons dissociate from the atoms and recombine with them, emitting light: in this case luminescence and incandescence become indistinguishable.

Luminescent pigments and dyes. Nonluminescent pigments and dyes exhibit colours because they absorb white light and reflect that part of the spectrum that is complementary to the absorbed light. A small fraction of the absorbed light is transformed into heat, but no appreciable radiation is produced. If, however, an appropriate luminescent pigment absorbs daylight in a special region of its spectrum, it can emit light of a colour different from that of the reflected light. This is the result of electronic processes within the molecule of the dye or pigment by which even ultraviolet light can be transformed to visible—e.g., blue—light. These pigments are used in such diverse ways as in outdoor advertising, blacklight displays, and laundering: in the latter case, a residue of the “brightener” is left in the cloth, not only to reflect white light but also to convert ultraviolet light into blue light, thus offsetting any yellowness and reinforcing the white appearance.

EARLY INVESTIGATIONS

Although lightning, the aurora borealis, and the dim light of glowworms and of fungi have always been known to mankind, the first investigations (1603) of luminescence began with a synthetic material, when Vincenzo Cascariolo, an alchemist and cobbler in Bologna, Italy, heated a mixture of barium sulfate (in the form of barite, heavy spar) and coal; the powder obtained after cooling exhibited a bluish glow at night, and Cascariolo observed that this glow could be restored by exposure of the powder to sunlight. The name *lapis solaris*, or “sunstone,” was given to the material because alchemists at first hoped it would transform baser metals into gold, the symbol for gold being the Sun. The pronounced afterglow aroused the interest of many learned men of that period, who gave the material other names, including phosphorus, meaning “light bearer,” which thereafter was applied to any material that glowed in the dark.

Today, the name phosphorus is used for the chemical element only, whereas certain microcrystalline luminescent materials are called phosphors. Cascariolo’s phosphor evidently was a barium sulfide; the first commercially available phosphor (1870) was “Balmain’s paint,” a calcium sulfide preparation. In 1866 the first stable zinc sulfide phosphor was described. It is one of the most important phosphors in modern technology.

One of the first scientific investigations of the luminescence exhibited by rotting wood or flesh and by glowworms, known from antiquity, was performed in 1672 by Robert Boyle, an English scientist, who, although not aware of the biochemical origin of that light, nevertheless established some of the basic properties of bioluminescent systems: that the light is cold; that it can be inhibited by chemical agents such as alcohol, hydrochloric acid, and ammonia; and that the light emission is dependent on air (as later established, on oxygen).

In 1885–87 it was observed that crude extracts prepared from West Indian fireflies (*Pyrophorus*) and from the boring clam, *Pholas*, gave a light-emitting reaction when mixed together. One of the preparations was a cold-water extract containing a compound relatively unstable to heat, luciferase; the other was a hot-water extract containing a relatively heat-stable compound, luciferin. The luminescent reaction that occurred when solutions of luciferase and luciferin were mixed at room temperature suggested that all bioluminescent reactions are “luciferin–luciferase reactions.” In view of the complex nature of bioluminescent reactions, it is not astonishing that this simple concept of bioluminescence has had to be modified. Only a small number of bioluminescent systems have been investigated for their respective luciferin and the corresponding luciferase, the best known being the bioluminescence of

Limitations of quantum mechanics

Luminescence a quantum process

Phosphors

Bioluminescence

fireflies from the United States, a little crustacean living in the Japanese sea (*Cypridina hilgendorfi*), and decaying fish and flesh (bacterial bioluminescence). Although bioluminescent systems have not yet found practical applications, they are interesting because of their high luminescence efficiency.

The first efficient chemiluminescent materials were non-biological synthetic compounds such as luminol (with the formula 5-amino-2,3-dihydro-1,4-phthalazinedione). The strong blue chemiluminescence resulting from oxidation of this compound was first reported in 1928.

PHOSPHORESCENCE AND FLUORESCENCE

The name luminescence has been accepted for all light phenomena not caused solely by a rise of temperature, but the distinction between the terms phosphorescence and fluorescence is still open to discussion. With respect to organic molecules, the term phosphorescence means light emission caused by electronic transitions between levels of different multiplicity (explained more fully below), whereas the term fluorescence is used for light emission connected with electronic transitions between levels of like multiplicity. The situation is far more complicated in the case of inorganic phosphors.

The term phosphorescence was first used to describe the persistent luminescence (afterglow) of phosphors. The mechanism described above for the phosphorescence of excited organic molecules fits this picture in that it is also responsible for light persistence up to several seconds. Fluorescence, on the other hand, is an almost instantaneous effect, ending within about 10^{-8} second after excitation. The term fluorescence was coined in 1852, when it was experimentally demonstrated that certain substances absorb light of a narrow spectral region (*e.g.*, blue light) and instantaneously emit light in another spectral region not present in the incident light (*e.g.*, yellow light) and that this emission ceases at once when the irradiation of the material comes to an end. The name fluorescence was derived from the mineral fluor spar, which exhibits a violet, short-duration luminescence on irradiation by ultraviolet light.

LUMINESCENCE EXCITATION

Chemiluminescence and bioluminescence. Most of the energy liberated in chemical reactions, especially oxidation reactions, is in the form of heat. In some reactions, however, part of the energy is used to excite electrons to higher energy states, and, for fluorescent molecules, chemiluminescence results. Studies indicate that chemiluminescence is a universal phenomenon, although the light intensities observed are usually so small that sensitive detectors are necessary. There are, however, some compounds that exhibit brilliant chemiluminescence, the best known being luminol, which, when oxidized by hydrogen peroxide, can yield a strong blue or blue-greenish chemiluminescence. Other instances of strong chemiluminescences are lucigenin (an acridinium compound) and lophine (an imidazole derivative). In spite of the brilliance of their chemiluminescence, not all of these compounds are efficient in transforming chemical energy into light energy, because only about 1 percent or less of the reacting molecules emit light. During the 1960s, esters (organic compounds that are products of reactions between organic acids and alcohols) of oxalic acid were found that, when oxidized in nonaqueous solvents in the presence of highly fluorescent aromatic compounds, emit brilliant light with an efficiency up to 23 percent.

Bioluminescence is a special type of chemiluminescence catalyzed by enzymes. The light yield of such reactions can reach 100 percent, which means that almost without exception every molecule of the reacting luciferin is transformed into a radiating state. All of the bioluminescent reactions best known today are catalyzed oxidation reactions occurring in the presence of air.

Triboluminescence. When crystals of certain substances—*e.g.*, sugar—are crushed, luminescent sparkles are visible. Similar observations have been made with numerous organic and inorganic substances. Closely related are the faint blue luminescence observable when adhesive tapes are stripped from a roll, and the luminescence

exhibited when strontium bromate and some other salts are crystallized from hot solutions. In all of these cases, positive and negative electric charges are produced by the mechanical separation of surfaces and during the crystallization process. Light emission then occurs by discharge, either directly, by molecule fragments, or via excitation of the atmosphere in the neighbourhood of the separated surface: the blue glow coming from adhesive tapes being unrolled is emitted from nitrogen molecules of the air that have been excited by the electric discharge.

Thermoluminescence. Thermoluminescence means not temperature radiation but enhancement of the light emission of materials already excited electronically by the application of heat. The phenomenon is observed with some minerals and, above all, with crystal phosphors after they have been excited by light.

Photoluminescence. Photoluminescence, which occurs by virtue of electromagnetic radiation falling on matter, may range from visible light through ultraviolet, X-ray, and gamma radiation. It has been shown that, in luminescence caused by light, the wavelength of emitted light generally is equal to or longer than that of the exciting light (*i.e.*, of equal or less energy). As explained below, this difference in wavelength is caused by a transformation of the exciting light, to a greater or lesser extent, to nonradiating vibration energy of atoms or ions. In rare instances—*e.g.*, when intense irradiation by laser beams is used or when sufficient thermal energy contributes to the electron excitation process—the emitted light can be of shorter wavelength than the exciting light (anti-Stokes radiation).

The fact that photoluminescence can also be excited by ultraviolet radiation was first observed by a German physicist, Johann Wilhelm Ritter (1801), who investigated the behaviour of phosphors in light of various colours. He found that phosphors luminesce brightly in the invisible region beyond violet and thus discovered ultraviolet radiation. The transformation of ultraviolet light to visible light has much practical importance.

Gamma rays and X rays excite crystal phosphors and other materials to luminescence by the ionization process (*i.e.*, the detachment of electrons from atoms), followed by a recombination of electrons and ions to produce visible light. Advantage of this is taken in the fluoroscope used in X-ray diagnostics and in the scintillation counter that detects and measures gamma rays directed onto a phosphor disk that is in optical contact with the face of a photomultiplier tube (a device that amplifies light signals).

Electroluminescence. Like thermoluminescence, the term electroluminescence includes several distinct phenomena, a common feature of which is that light is emitted by an electrical discharge in gases, liquids, and solid materials. Benjamin Franklin, in the United States, for example, in 1752 identified the luminescence of lightning as caused by electric discharge through the atmosphere. An electric-discharge lamp was first demonstrated in 1860 to the Royal Society of London. It produced a brilliant white light by the discharge of high voltage through carbon dioxide at low pressure. Modern fluorescent lamps are based on a combination of electroluminescence and photoluminescence: mercury atoms in the lamp are excited by electric discharge, and the ultraviolet light emitted by the mercury atoms is transformed into visible light by a phosphor.

The electroluminescence sometimes observed at the electrodes during electrolysis is caused by the recombination of ions (therefore, this is a sort of chemiluminescence). The application of an electric field to thin layers of luminescing zinc sulfide can produce light emission, which is also called electroluminescence.

A great number of materials luminesce under the impact of accelerated electrons (once called cathode rays)—*e.g.*, diamond, ruby, crystal phosphors, and certain complex salts of platinum. The first practical application of cathodoluminescence was in the viewing screen of an oscilloscope tube constructed in 1897; similar screens, employing improved crystal phosphors, are used in television, radar, oscilloscopes, and electron microscopes.

The impact of accelerated electrons on molecules can

Role of
enzymes

Electric
discharge

produce molecular ions, ions of molecule fragments, and atomic ions. In gas-discharge tubes these particles were first detected as "canal rays" or anode rays. They are able to excite phosphors but not as efficiently as electrons can.

Radioluminescence. Radioactive elements can emit alpha particles (helium nuclei), electrons, and gamma rays (high-energy electromagnetic radiation). The term radioluminescence, therefore, means that an appropriate material is excited to luminescence by a radioactive substance. When alpha particles bombard a crystal phosphor, tiny scintillations are visible to microscopic observation. This is the principle of the device used by an English physicist, Ernest Rutherford, to prove that an atom has a central nucleus. Self-luminous paints, such as are used for dial markings for watches and other instruments, owe their behaviour to radioluminescence. These paints consist of a phosphor and a radioactive substance, *e.g.*, tritium or radium. An impressive natural radioluminescence is the aurora borealis: by the radioactive processes of the sun, enormous masses of electrons and ions are emitted into space in the solar wind. When they approach the Earth, they are concentrated by its geomagnetic field near the poles. Discharge processes of the particles in the upper atmosphere yield the famous luminance of the auroras.

LUMINESCENT MATERIALS AND PHOSPHOR CHEMISTRY

The first phosphor synthesized was probably an impure barium sulfide preparation with very low luminance efficiency and with the serious shortcoming that it was rather quickly decomposed in moist air, yielding hydrogen sulfide. A more stable sulfide-type phosphor was produced in 1866 by heating zinc oxide in a stream of hydrogen sulfide. In 1887 it became known that these sulfides do not luminesce in a chemically pure state but only when they contain small quantities of a so-called activator metal. Later, other materials, such as certain metal oxides, silicates, and phosphates, were found to luminesce if they were prepared by special procedures.

Sulfide-type phosphors, activators, fluxes. The sulfides of zinc and of cadmium are the most important basic materials of sulfide-type phosphors. An important condition of getting highly efficient phosphors is that these sulfides must first be prepared to the highest possible chemical purity before the necessary amount of activator can be added precisely. The emission of zinc sulfide can be shifted to longer wavelengths by increasing substitution of the zinc ions by cadmium ions. Zinc sulfide and cadmium sulfide phosphors are especially efficient in electroluminescence.

Sulfide-type phosphors are produced from pure zinc or cadmium sulfide or their mixtures by heating them together with small quantities (0.1–0.001 percent) of copper, silver, gallium, or other salts (activators) and with about 2 percent of sodium or another alkali chloride at about 1,000° C (1,832° F). The role of the alkali halides is to facilitate the melting process and, above all, to serve as coactivators (fluxes). Only small quantities of the alkali halide are integrated into the phosphor, but this small quantity is highly important for its luminescence efficiency. Copper-activated zinc and cadmium sulfides exhibit a rather long afterglow when their irradiation has ceased, and this is favourable for application in radar screens and self-luminous phosphors.

Oxide-type phosphors. Certain oxide-type minerals have been found to luminesce when irradiated. In some of them, activators must first be introduced into the crystal. Examples are ruby (aluminum oxide with chromium activator—bright-red emission) and willemite (zinc orthosilicate with manganese activator—green emission). On the other hand, scheelite (calcium tungstate) emits a blue luminescence without activator. All of these minerals have been made synthetically, with remarkably higher efficiencies than those that occur naturally. Silicates, borates, and phosphates of the second group of the periodic table of elements, such as zinc silicate, zinc beryllium silicate, zinc and cadmium borates, and cadmium phosphates, become efficient phosphors when activated with manganese ions, emitting in the red to green region of the spectrum. They have been incorporated into colour television screens to emit the colours blue (silver-activated zinc sul-

fide), green (manganese-activated zinc orthosilicate), and red (europium-activated yttrium vanadate).

Centres, activators, coactivators, poisons. The study of phosphor chemistry has yielded a detailed picture of the role of the above-mentioned activators and fluxes. Philipp Anton Lenard, a physicist in Germany, was the first (1890) to describe activator ions as being distributed in zinc sulfide and other crystalline materials that serve as the host crystal. The activator ions are surrounded by host-crystal ions and form luminescing centres where the excitation–emission process of the phosphor takes place. These centres must not be too close together within the host crystal lest they inactivate each other. For high efficiency, only a trace of the activator may be inserted into the host crystal, and its distribution must be as regular as possible. In high concentration, activators act as "poisons" or "killers" and thus inhibit luminescence. The term killer is used especially for iron, cobalt, and nickel ions, whose presence, even in small quantities, can inhibit the emission of light from phosphors.

Killers

Phosphors, such as calcium tungstate or zinc sulfide, that need no activator appear to have their luminescing centres in special groups of atoms different from the symmetry of their own crystal lattice, such as the group WO_4 in the compound calcium tungstate ($CaWO_4$), or, similarly, the SiO_4 group in zinc orthosilicate, (Zn_2SiO_4). That luminescing properties of a centre are strongly dependent on the symmetry of neighbouring ion groups with respect to the whole phosphor molecule is clearly proved by the spectral shifts of certain phosphors activated with lanthanide ions, which emit in narrow spectral regions. Because of this altering effect on the symmetry of luminescing centres, small quantities (about 0.2 percent) of titania incorporated in zinc orthosilicate give a remarkable increase in luminescence. Titania is called an intensifier activator because it increases the host-crystal luminescence, whereas a substance that produces luminescence not exhibited by the chemically pure host crystal is called an originative activator.

The fluxes (*e.g.*, sodium chloride) act as coactivators by facilitating the incorporation of activator ions. Copper ions, for instance, are used as activators of zinc chloride phosphors and are usually introduced in the copper(II), or cupric, form (the Roman numeral indicates the oxidation state; that is, I means that the element has one electron involved in a chemical bond and II that it has two electrons involved; the larger oxidation state is indicated by the *-ic* ending and the smaller by the *-ous* ending). If a copper(II) compound is incorporated into the zinc sulfide by heating, copper(I) sulfide (or cuprous sulfide, formula Cu_2S) will be produced with crystals that will not fit into the host-crystal zinc chloride because their form is so different, and only a relatively few luminescent centres will be possible. On the other hand, if a coactivator such as sodium chloride is introduced along with the copper(II) salt, the copper(II) ions are reduced to form copper(I) chloride (or cuprous chloride, formula $CuCl$) crystals with the same structure as the host crystal. Thus, many luminescent centres will be produced, and strong activation will result.

In describing a luminescent phosphor, the following information is pertinent: crystal class and chemical composition of the host crystal, activator (type and percentage), coactivator (intensifier activator), temperature and time of crystallization process, emission spectrum (or at least visual colour), and persistence. A few phosphors and their activators are listed in Table 2.

Organic luminescent materials. Although the inorganic phosphors are industrially produced in far higher quantities (several hundred tons per year) than the organic luminescent materials, some types of the latter are becoming more and more important in special fields of practical application. Paints and dyes for outdoor advertising contain strongly fluorescing organic molecules such as fluorescein, eosin, rhodamine, and stilbene derivatives. Their main shortcoming is their relatively poor stability in light, because of which they are used mostly when durability is not required. Organic phosphors are used as optical brighteners for invisible markers of laundry, banknotes, identity cards, and stamps and for fluorescence microscopy of tis-

Table 2: Visual Properties of Some Luminescent Materials

phosphor (phosphor/activator; coactivator)	emission	
	colour*	persistence
rhomboidal Zn ₂ SiO ₄ /Mn 0.3% (rhomboidal zinc orthosilicate/ manganese 0.3%); 1,200° C 60 min, slow cooling	green (525 nm)	short (0.01 sec)
β-Zn ₂ SiO ₄ /Mn 0.3% (beta zinc orthosilicate/manganese 0.3%); 1,600° C 10 min, quench cooling	yellow (610 nm)	short (0.01 sec)
cubic ZnS/Cu 0.03%; Cl (cubic zinc sulfide/copper 0.03%; chloride); 950° C 10 min, slow cooling	green-blue (516 nm)	long (hours)
hexagonal ZnS/Cu 0.03%; Cl (hexagonal zinc sulfide/ copper 0.03%; chloride); 1,250° C 10 min, slow cooling	green (528 nm)	very long (up to 24 hours)

*The wavelengths of the respective emission maxima are given in parentheses. 1 nanometre (nm) = 10⁻⁹ metre = 1 millimicron = 10 angstroms.

sues in biology and medicine. Their "invisibility" is due to the fact that they absorb practically no visible light. The fluorescence is excited by invisible ultraviolet radiation (black light).

Photoradiation in gases, liquids, and crystals. When describing chemical principles associated with luminescence, it is useful, at first, to neglect interactions between the luminescing atoms, molecules, or centres with their environment. In the gas phase these interactions are smaller than they are in the condensed phase of a liquid or a solid material. The efficiency of luminescence in the gas phase will be far greater than in the condensed phases because in the latter the energy of the electrons excited by photons or by chemical-reaction energy can be dissipated as thermal, nonradiative energy by collision of the atoms or by the rotational and vibrational energy of the molecules. This effect has to be taken into account even more when the radiation of single atoms is compared with that of multi-atomic molecules. For molecules, radiative (electronic-excitation) energy is internally converted to vibrational energy; that is, there are radiationless transitions of electrons in atoms. This is the explanation for the fact that only a relatively small number of compounds are able to exhibit efficient luminescence. In crystals, on the other hand, the binding forces between the ions or atoms of the lattice are strong compared with the forces acting between the particles of a liquid, and electron-excitation energy, therefore, is not as easily transformed into vibrational energy, thus leading to a good efficiency for radiative processes.

LUMINESCENCE PHYSICS

Mechanism of luminescence. The emission of visible light (that is, light of wavelengths between about 690 nanometres and 400 nanometres, corresponding to the region between deep red and deep violet) requires excitation energies the minimum of which is given by Einstein's law stating that the energy (E) is equal to Planck's constant (h) times the frequency of light (ν), or Planck's constant times the velocity of light (c) in a vacuum divided by its wavelength (λ); that is,

$$E = h\nu = \frac{hc}{\lambda}$$

The energy required for excitation therefore ranges between 40 kilocalories (for red light), about 60 kilocalories (for yellow light), and about 80 kilocalories (for violet light) per mole of substance. Instead of expressing these energies in kilocalories, electron volt units (one electron volt = 1.6×10^{-12} erg; the erg is an extremely small unit of energy) may be used, and the photon energy thus required in the visible region ranges from 1.8 to 3.1 electron volts.

The excitation energy is transferred to the electrons responsible for luminescence, which jump from their ground-state energy level to a level of higher energy. The energy levels that electrons can assume are specified by

quantum mechanical laws. The different excitation mechanisms considered below depend on whether or not the excitation of electrons occurs in single atoms, in single molecules, in combinations of molecules, or in a crystal. They are initiated by the means of excitation described above: impact of accelerated particles such as electrons, positive ions, or photons. Often, the excitation energies are considerably higher than those necessary to lift electrons to a radiative level; for example, the luminescence produced by the phosphor crystals in television screens is excited by cathode-ray electrons with average energies of 25,000 electron volts. Nevertheless, the colour of the luminescent light is nearly independent of the energy of the exciting particles, depending chiefly on the excited-state energy level of the crystal centres.

Electrons taking part in the luminescence process are the outermost electrons of atoms or molecules. In fluorescent lamps, for example, a mercury atom is excited by the impact of an electron having an energy of 6.7 electron volts or more, raising one of the two outermost electrons of the mercury atom in the ground state to a higher level. Upon the electron's return to the ground state, an energy difference is emitted as ultraviolet light of a wavelength of 185 nanometres. A radiative transition between another excited state and the ground-state level of the mercury atom produces the important ultraviolet emission of 254-nanometre wavelength, which, in turn, can excite other phosphors to emit visible light. (One such phosphor frequently used is a calcium halophosphate incorporating a heavy-metal activator.)

This 254-nanometre mercury radiation is particularly intensive at low mercury vapour pressures (around 10^{-5} atmosphere) used in low-pressure discharge lamps. About 60 percent of the input electron energy may thus be transformed into near-monochromatic ultraviolet light; *i.e.*, ultraviolet light of practically one single wavelength.

Whereas at low pressure there are relatively few collisions of mercury atoms with each other, the collision frequency increases enormously if mercury gas is excited under high pressure (*e.g.*, eight atmospheres or more). Such excitation leads not only to collisional de-excitation of excited atoms but also to additional excitation of excited atoms. As a consequence, the spectrum of the emitted radiation no longer consists of practically one single, sharp spectral line at 254 nanometres, but the radiation energy is distributed over various broadened spectral lines corresponding to different electronic energy levels of the mercury atom, the strongest emissions lying at 303, 313, 334, 366, 405, 436, 546, and 578 nanometres. High-pressure mercury lamps can be used for illumination purposes because the emissions from 405 to 546 nanometres are visible light of bluish-green colour; by transforming a part of the mercury line emission to red light by means of a phosphor, white light is obtained.

When gaseous molecules are excited, their luminescence spectra show broad bands; not only are electrons lifted to levels of higher energy but vibrational and rotational motions of the atoms as a whole are excited simultaneously. This is because vibrational and rotational energies of molecules are only about 10^{-2} and 10^{-4} , respectively, those of the electronic transition energies, and these many energies can be added to the energy of a single electronic transition, which is represented by a multitude of slightly different wavelengths making up one band. In larger molecules, several overlapping bands, one for each kind of electronic transition, can be emitted. Emission from molecules in solution is predominantly bandlike caused by interactions of a relatively great number of excited molecules with molecules of the solvent. In molecules, as in atoms, the excited electrons generally are outermost electrons of the molecular orbitals.

The terms fluorescence and phosphorescence can be used here, on the basis not only of the persistence of luminescence but also of the way in which the luminescence is produced. When an electron is excited to what is called, in spectroscopy, an excited singlet state, the state will have a lifetime of about 10^{-8} second, from which the excited electron can easily return to its ground state (which normally is a singlet state, too), emitting its excitation energy

Effect of
binding
force

Electron
excitation

as fluorescence. During this electronic transition the spin of the electron is not altered; the singlet ground state and the excited singlet state have like multiplicity (number of subdivisions into which a level can be split). An electron, however, may also be lifted, under reversal of its spin, to a higher energy level, called an excited triplet state. Singlet ground states and excited triplet states are levels of different multiplicity. For quantum mechanical reasons, transitions from triplet states to singlet states are "forbidden," and, therefore, the lifetime of triplet states is considerably longer than that of singlet states. This means that luminescence originating in triplet states has a far longer duration than that originating in singlet states; phosphorescence is observed.

Phosphorescence

The interactions of a large number of atoms, ions, or molecules are greater still in solution and in solids; to obtain a narrowing of the spectral band, subzero temperatures (down to that of liquid helium) are applied in order to reduce vibrational motions. The electronic energy levels of crystals such as zinc sulfide and other host crystals used in phosphors form bands; in the ground state practically all electrons are to be found on the valence band, whereas they reach the conduction band after sufficient excitation (see MATTER). The energy difference between the valence band and the conduction band corresponds to photons in the ultraviolet or still shorter wavelength region. Additional energy levels are introduced by activator ions or centres bridging the energy gap between valence band and conduction band, and, when an electron is transferred from the valence band to such an additional energy level by excitation energy, it can produce visible light on return to the ground state. A rather close analogy exists between the forbidden transitions of certain excited molecular electronic states (triplet-singlet, leading to phosphorescence) and the transition of an electron of an inorganic phosphor kept in a trap: traps (certain distortions in the crystal lattice) are places in the crystal lattice where the energy level is lower than that of the conduction band, and from which the direct return of an electron to the ground state is also forbidden.

When a solid is bombarded by photons or particles, the excitation of the centres can occur directly or by energy transfer. In the latter case, excited but nonluminescing states are produced at some distance from the centre, with the energy moving through the crystal in the form of excitons (ion-electron pairs) until it approaches a centre where the excitation process can occur. This energy transfer can also be realized by radiation in inorganic phosphors containing two activators, as well as in solutions of organic molecules.

Spontaneous and stimulated emission. The radiative return of excited electrons to their ground state occurs spontaneously, and when there exists an assembly of excited electrons their individual spontaneous radiative transitions are independent of each other. Therefore, the luminescence light is incoherent (the emitted waves are not in phase with each other) in this case. Sometimes the emission of luminescence can be stimulated by irradiation with photons of the same frequency as that of the emitted light; such stimulated transitions are used in lasers, which produce very intensive beams of coherent monochromatic light.

The spontaneous luminescent emission follows an exponential law that expresses the rate of intensity decay and is similar to the equation for the decay of radioactivity and some chemical reactions. It states that the intensity of luminescent emission is equal to an exponential value of minus the time of decay divided by the decay time, or $L = L_0 \exp(-t/\tau)$, in which L is the intensity of emission at a time t after an initial intensity L_0 , and τ is the decay time of the luminescence; that is, the time in which the assembly of the excited atoms would decrease in luminescence intensity to a value of $0.368 L_0$.

When excited atoms of the centres are in contact with other atoms, as is the case in condensed phases (liquids, solids, in gases of not-too-low pressure), part of the excitation energy will be transformed into heat by collisional deactivation (thermal quenching). The decay time, therefore, has to be replaced by an effective excited-state lifetime, re-

sulting in a more complicated exponential decay law that depends on the collision frequency, the energy imparted to the excited atoms of the centre that causes the transfer of excitation energy into heat (activation energy), a constant, and the temperature of the luminescent material. This law describes the actual luminescence decay of a great number of luminescent materials; e.g., calcium tungstate.

Increase of activation energy for nonradiative deactivation of excited-centres luminescence decay can be achieved by changing the host crystal or by electron traps. The traps are imperfections in the crystal lattice where electrons are captured after they have been ejected from a luminescent centre by excitation energy. That the luminescent properties of phosphor centres are strongly dependent on the chemical nature of the host crystal may be seen in Table 3, showing that the same activator ions (manganese ions with two positive charges, indicated as Mn^{2+} , or $Mn[II]$), in different host crystals yield remarkably different-coloured emissions and decay times (measured in fractions of a second).

Prolonging the emission time of phosphors up to days or even longer (production of phosphorescence of the phosphors) is possible by inserting traps into the host crystal. Trapped electrons cannot return directly to the centre. In order to be released from the traps they must first obtain additional thermal energy—in this case, thermal energy stimulates luminescence—after which they recombine with a centre and undergo radiative transition. Trapping in crystals has its analogy to forbidden transitions in

Table 3: Influence of Host Crystal on the Lifetime and Emission Colour of the Excited Phosphors

host crystal	activator	time (second)	emission colour
Tetragonal zinc fluoride, ZnF_2	manganese(II)	0.1	orange
Rhombic cadmium sulfate, $CdSO_4$	manganese(II)	0.05	orange
Rhombic magnesium sulfate, $MgSO_4$	manganese(II)	0.03	red
Rhombic zinc phosphate, $Zn_3(PO_4)_2$	manganese(II)	0.02	red
Cadmium silicate, $CdSiO_3$	manganese(II)	0.019	orange
Zinc orthosilicate, Zn_2SiO_4	manganese(II)	0.018	yellow
Cadmium pyroborate, $Cd_2B_2O_5$	manganese(II)	0.015	red orange
Rhombohedral zinc orthosilicate, Zn_3SiO_4	manganese(II)	0.013	green
Rhombohedral zinc germanate, Zn_3GeO_4	manganese(II)	0.0105	green yellow
Cubic zinc aluminate, $ZnAl_2O_4$	manganese(II)	0.0055	blue green
Cubic zinc gallate, $ZnGa_2O_4$	manganese(II)	0.0043	green blue
Hexagonal zinc sulfide, ZnS	manganese(II)	0.0004	orange

molecules (triplet-singlet transitions) or in radiation processes from metastable atomic-energy levels.

An example of a practical application of stimulated emission of a phosphor with trapped electrons is cubic strontium sulfide/selenide activated with samarium and europium ions, the coactivators being strontium sulfate and calcium fluoride. This phosphor has been used in devices for viewing scenes at night by reflected infrared light emitted by infrared lamps. The traps in this phosphor have been identified as samarium ions, whereas europium ions are the active ions in the centres. The phosphor is first excited by photons of about three electron volts (blue light), which results in an ejection of an electron from a europium ion (Eu^{2+}) centre. This excited electron is trapped by a triply charged samarium ion (Sm^{3+}), which is transferred to a doubly charged samarium ion (Sm^{2+}). Heat or irradiation by infrared photons releases one electron from the doubly charged samarium ion (Sm^{2+}). The electron is then recaptured from a triply charged europium ion (Eu^{3+}), yielding an excited doubly charged europium ion (Eu^{2+}), which returns to its ground state by emitting a photon of 2.2 electron volts energy (yellow light). The trap depth of this phosphor (*i.e.*, the energy required for release of an electron from it) is large compared to the thermal

Exponential decay

energy of the lattice of the host crystal, and, therefore, the lifetime of the traps at room temperature is many months long. Bombarding this phosphor with photons of energy higher than that of infrared photons but not sufficient for excitation can lead to photoquenching: the traps are emptied far more rapidly, and thermal deactivation of the centres is enhanced.

When iron, cobalt, or nickel ions are present in a phosphor, an excited electron can be captured by these ions. The excitation energy is then emitted as infrared photons, not as visible light, so that luminescence is quenched. These ions, therefore, are called killers—the killing process being opposite to stimulation.

In chemiluminescence, such as the oxidation of luminol, light emission depends not only on radiative and quenching or intramolecular deactivation processes but also on the efficiency of the chemical reaction leading to molecules in an electronically excited state.

In bioluminescence reactions, the production of electronically excited molecules, as well as their radiative transitions back to their ground state, is efficiently catalyzed by the enzymes acting here, and bioluminescence light output is therefore high.

The luminescence photons emitted by one kind of excited atom, molecule, or phosphor can excite another to emit its specific luminescence: this type of energy transfer is observed with inorganic as well as organic substances. Thus, excited benzene molecules can excite naphthalene molecules by radiative-energy transfer. The radiation produced by the luminol chemiluminescence can produce fluorescence when fluorescein is added to the reaction mixture. In most of these cases the acceptor molecules have luminescent electrons with energy levels lower than those of the primary excited molecules, and emitted secondary luminescence is therefore of longer wavelength than the primary. Practical application of this phenomenon, called cascading, is used in radar kinescopes, which have composite fluorescent screens consisting of a layer of blue-emitting zinc sulfide/silver (chloride) phosphor—the hexagonal crystal, ZnS/Ag(Cl) deposited on a layer of yellow-emitting zinc or cadmium sulfide/copper [chloride] phosphor [the hexagonal crystal, (Zn,Cd)S/Cu (Cl)].

The cathode-ray electrons excite the blue-emitting phosphor, whose photons, in turn, excite the yellow-emitting phosphor, which has traps with a decay time of about 10 seconds. Excitation of the blue-emitting phosphor alone would be unfavourable, as the sharply focussed cathode rays are absorbed by the blue phosphor to a small extent only, and its decay time is too short; also, direct excitation of the yellow-emitting phosphor alone would yield poor efficiency because the traps are emptied too rapidly by the heat produced by the relatively high-energetic electron impact.

Another energy-transfer mechanism is referred to as sensitization: a calcium carbonate phosphor (rhombohedral CaCO_3/Mn), for example, emits orange light under cathode-ray irradiation but is not excited by the 254-nanometre emission of mercury atoms, whereas this emission produces the same orange light with calcium carbonate (rhombohedral CaCO_3) activated by manganese and lead ions. This is not cascade luminescence: a mechanical mixture of a manganese and a lead-activated calcium carbonate exhibits no emission under ultraviolet radiation. In a phosphor containing both activators, the lead ions act as sensitizers in introducing an additional excitation band into the system from which the manganese ions get their excitation energy in a nonradiative energy transfer. Similar sensitization is observed in gases and in liquids.

Solid state energy states. The complicated problems concerning the energy states in solids of a luminescence centre are commonly visualized by adapting the energy-level diagram used in describing energy transitions in an isolated diatomic molecule (Figure 27).

In this diagram, the potential energy of a centre is plotted as a function of the average distance (\bar{x}) between the atoms: \bar{x}^* represents the ground state and \bar{x}_0 represents the lowest excited state of the centre. In a tetrahedral permanganate-ion centre (MnO_4), for example \bar{x} would be the average distance between the central manganese ion

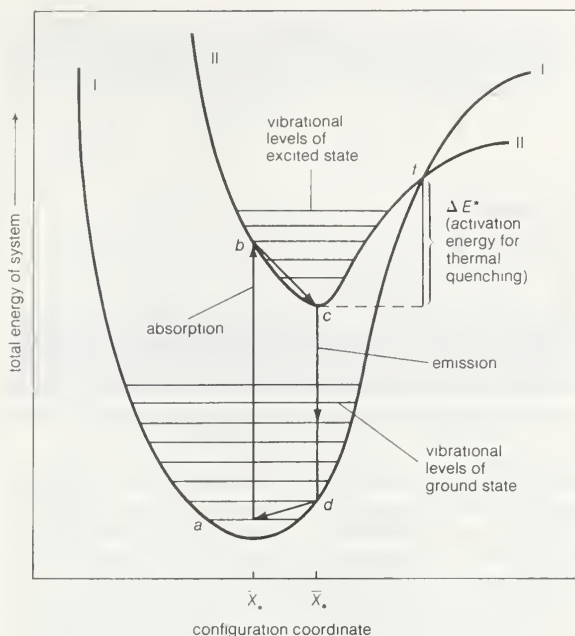


Figure 27: Energy levels of a luminescent centre (see text).

and an oxygen ion in any of the corners of the tetrahedron.

At a temperature of absolute zero the ground-state energy level is near the bottom of curve I at the minimum amplitude of atomic vibration. At room temperature (300 K [81° F]) the ground state lies higher, at *a*, where the centre has considerable vibrational energy. When an electron of the centre is excited, it is lifted to the higher energy level at *b* in curve II. This electronic transition occurs far more rapidly than the readjustment of the atoms of the centre, which then occurs within a time of about 10^{-12} second to reach the minimum vibrational level at *c*. The energy difference (*b* – *c*) is dissipated as heat in the host-crystal lattice. From the excited-state level *c*, the electron can return to the ground-state level *d* shown in Curve I, the liberated energy being emitted as a photon.

The last step is a readjustment of the centre to *a*, the energy difference (*d* – *a*) again being dissipated as heat. Nonradiative transition of the excited electron back to its ground state occurs when the electron is excited to an energy level above the intersection point *f* of the ground-state and the excited-state energy curve. This is caused mainly by increasing the vibrations of the lattice by application of higher temperatures. The energy difference (*f* – *c*) is equal to the activation energy already mentioned, and therefore most centres become increasingly nonradiative at higher temperatures. In trap-type phosphors the temperature must be sufficiently high, of course, to eject the electron from the traps.

In some phosphors—calcium tungstate (CaWO_4), for example—absorption and emission of the exciting energy appear to take place mainly in the same centre; the excited electron remains near the centre. Such phosphors do not exhibit photoconductivity because only a few excited electrons succeed in reaching the conduction band where they are freely mobile. The luminescence decay is exponential.

Zinc sulfide phosphors, however, are photoconducting, which means that many excited electrons are lifted to the conduction band of the host crystal. The energy levels of different centres and of the host-crystal lattice have to be taken into account simultaneously.

The relative levels of the zinc sulfide valence band (ground state of the host-crystal lattice) and the conduction band (excited state of the host-crystal lattice), of activator levels and of trap levels are shown in Figure 28. Points 1, 2, 3, and 4 represent one situation in a host crystal, and points 5, 6, 7, 8, 9, and 10 represent another situation.

The activator ions introduce additional ground-state levels and excited-state levels of energies between those of the valence and the conduction band of the zinc sulfide. When the excitation energy is sufficiently high, an electron is raised to the conduction band (1 → 2, 5 → 6, corre-

Cascading

Photocon-
ducting
phosphors

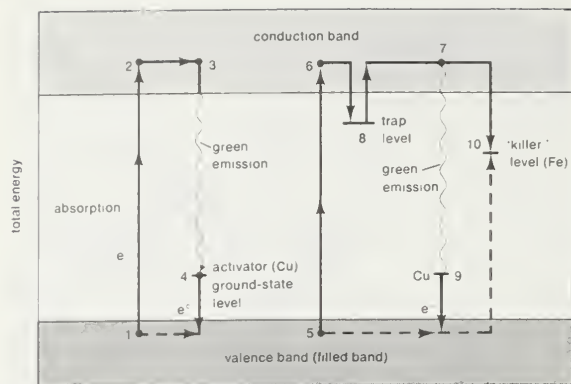


Figure 28: Transition of an electron from the valence band to the conduction band by light absorption (see text).

sponding to the ionization continuum in a gas). It moves away from the centre (2 → 3, 6 → 8) and may either be trapped by an imperfection of the lattice (8) or return to an ionized centre (activator), in which it first occupies an excited level (3 → 4) and then drops to the ground state of the activator centre by emitting a photon. An activator centre that captures such an excited electron has already lost one of its own electrons to a positive hole (electron vacancy) in the host-crystal lattice.

The energetic level of the traps is about 0.25 electron volt beneath the conduction-band level. A trapped electron (8) must be raised to the conduction band by thermal energy before a recombination with an ionized activator centre can occur. The green emission (530 nanometres) of the zinc sulfide phosphor (ZnS/Cu) is explained by the recombination of an electron from the conduction band and a copper ion in an activator centre (7 → 9); the blue emission (463 nanometres) is due to recombination of the excited electron and a copper ion in an interstitial place.

Direct excitation of the activator centres is also possible. When an electron recombines with a killer ion (10), no visible emission occurs.

In solid-state electroluminescence, the radiative processes occurring in a phosphor under irradiation are produced by applying external electric fields of several hundred volts, alternating at several thousand cycles per second. Special preparations of zinc sulfide (hexagonal ZnS), with an iodine coactivator and high concentrations of a copper activator, are embedded in a thin layer of about 0.01 centimetre (0.004 inch) of insulating organic material or glass, which is mounted between the electrodes.

High luminescence efficiencies result. Application of a direct-current field yields luminescence in crystals of gallium arsenide (GaAs), silicon carbide (SiC), cadmium sulfide (CdS), and zinc monocrystals of sulfide with copper activator (ZnS/Cu); the cathode injects electrons into the conduction band, whereas the anode removes electrons.

Efficiency of luminescence; luminance. The efficiency of luminescence emission must be regarded on an energy and a quantum basis. When every exciting photon yields an emitted photon of the same energy (as is the case for resonance excitation—*i.e.*, excitation of fluorescence by a monochromatic light of exactly the same wavelengths as the resulting fluorescence—and radiation of isolated atoms in dilute gases), the luminescence efficiency is 100 percent with respect to input energy as well as to the number of quanta. When the number of secondary photons is equal to that of the primary but their energy is less because some energy is dissipated as heat, the quantum efficiency is 100 percent but the luminescence efficiency is less than 100 percent. The quantum efficiency of most luminescences is far lower than 100 percent; zinc sulfide phosphors have about 20 percent efficiency, and solid-state electroluminescence is less than 10 percent efficient.

In chemiluminescence the quantum efficiency is about 1

percent in "brilliant" reactions, such as the oxidation of luminol, and up to 23 percent in the oxalate chemiluminescence. Solid-state electroluminescence, or electroluminescence of gases excited by high-frequency electric fields, is usually less than 10 percent.

The light intensity of luminescent processes depends chiefly on the excitation intensity, the density, and the lifetime of the radiative atoms, molecules, or centres. For practical purposes this luminous intensity per unit area is called photometric brightness or luminance of a material and is measured in lambert or millilambert (0.001 lambert) units (one lambert is equal to one candle per square centimetre divided by π).

(K.-D.G.)

BIBLIOGRAPHY

General works. A.C.S. VAN HEEL and C.H.F. VELZEL, *What is Light?* (1968; originally published in Dutch, 1968), is a good book for the general reader; numerous illustrations are accompanied by a simple text. MARCEL MINNAERT, *Light and Colour in the Outdoors* (1993; originally published in Dutch, 1937), also recommended for the general reader, deals with optical effects such as rainbows, mirages, and many other effects of light and colour. MAX BORN and EMIL WOLF, *Principles of Optics*, 6th ed. (1980, reprinted 1993), is a comprehensive treatise on the wave theory of light, particularly strong on diffraction, scattering, polarization, and coherence—suitable for the reader with advanced mathematical knowledge. R.W. DITCHBURN, *Light* (1953, reprinted 1991), is a text at the university level that includes both theory and experimental evidence in relation to wave phenomena and quantum optics. K.D. FROOME and L. ESSEN, *The Velocity of Light and Radio Waves* (1969), critically reviews modern methods of measuring the constant *c*. The work by J.H. SANDERS, *Velocity of Light* (1965), contains reprints of original papers by A.A. Michelson, F.G. Pease, F. Pearson, Louis Essen, A.C. Gordon-Smith, and E. Bergstrand, together with a review of other work on the subject. The history of early theories is traced in A.I. SABRA, *Theories of Light, From Descartes to Newton* (1967, reissued 1981). RICHARD MORRIS, *Light* (1979), provides a general introduction with good bibliographies. MICHAEL I. SOBEL, *Light* (1987), is an interesting, understandable, yet technical account with historical and practical examples.

Quantum mechanics and electromagnetic theory. LEONARD I. SCHIFF, *Quantum Mechanics*, 3rd ed. (1968), is a systematic exposition of the physical concepts of quantum mechanics, including a description of the electromagnetic field and its interactions with matter. JOSEPH NEEDHAM, *Science and Civilization in China*, vol. 4, pt. 1 (1962), includes an authoritative account of the study of light in ancient China and a detailed comparison with parallel developments due to Greek and Arab philosophers. EDMUND WHITTAKER, *A History of the Theories of Aether & Electricity*, rev. and enlarged ed., 2 vol. (1951–53, reissued in 1 vol., 1989), deals with the background to the development of the electromagnetic theory of light in the 19th century. F.H. READ, *Electromagnetic Radiation* (1980), includes lucid explanations of the basic concepts.

Luminescence. E. NEWTON HARVEY, *A History of Luminescence from the Earliest Times Until 1900* (1957), a classical work, deals rather extensively with the historical development of the different types of luminescence, especially bioluminescence and chemiluminescence. E.J. BOWEN (ed.), *Luminescence in Chemistry* (1968), is a textbook for students and researchers. C.A. PARKER, *Photoluminescence of Solutions* (1968), explains in detail the basic principles of luminescence as applied to photoluminescence in solutions, kinetics, apparatus, and analytical applications. M. ZANDER, *Phosphorimetry* (1968), is the first modern monograph dealing exclusively with the phosphorescence of organic materials, with a complete bibliography. GEORGE G. GUILBAULT (ed.), *Fluorescence: Theory, Instrumentation, and Practice* (1967), is a fairly technical account written by outstanding specialists in their respective fields (good background knowledge is necessary). DAVID M. HERCULES (ed.), *Fluorescence and Phosphorescence Analysis* (1966), comprises chapters of different grades of detail covering the luminescence field. G.F.J. GARLICK, "Luminescence," *Handbuch der Physik*, vol. 26, pp. 1–128 (1958), an extended text (written in English), systematically explains the physical phenomena and theory of luminescence. MARVIN C. GOLDBERG (ed.), *Luminescence Applications in Biological, Chemical, Environmental, and Hydrological Sciences* (1989), is a complete compilation of original research in a wide variety of areas. (R.W.D./K.-D.G./Ed.)

Lima

Lima, the national capital of the Republic of Peru, is also the country's commercial and industrial centre. Central Lima is located at an elevation of 512 feet (156 metres) on the south bank of the Río Rímac, about eight miles (13 kilometres) inland from the Pacific Ocean port of Callao, and has an area of 27 square miles (70 square kilometres). Its name is a corruption of the Quechua Indian name Rímac, meaning "Talker." The city forms a modern oasis, surrounded by the Peruvian coastal desert and overshadowed by the neighbouring Andes mountains.

This article is divided into the following sections:

Physical and human geography	29
The character of the city	29
The landscape	29
The city site	
Climate	
The city layout	
The people	30
The economy	30
Industry and commerce	
Transportation	
Administration and social conditions	31
Government	
Services	
Cultural life	31
History	32
Pre-Columbian and colonial periods	32
The modern city	32
Bibliography	32

Physical and human geography

THE CHARACTER OF THE CITY

Perhaps the best clue to the significance of Lima to the nation of Peru can be found in its most popular nickname: El Pulpo ("The Octopus"). Metropolitan Lima's huge size—it accounts for about one-third of the total population of Peru—has both resulted from and stimulated the concentration of people, capital, political influence, and social innovations. Lima's unique status is but one of the more important consequences of a highly centralized, unitary state that from its inception in the early 19th century solved interregional conflicts by focusing power and prestige on the city. With its port of Callao and its location at the centre of Peru's Pacific coast, Lima was long the only point of contact between the nation and the outside world.

As with many sprawling and rapidly growing metropolitan centres, Lima has its detractors as well as its promoters. Those who remember the more tranquil, traditional days, before the arrival of millions of migrants and before the many buses and automobiles brought pollution and congestion, are prone to use another nickname for the capital: Lima la Horrible. This is the noisy, dirty, gloomy, damp, and depressing Lima, perceptions shared by both short-term visitors and longtime residents. Even though sunshine does break through the dense coastal cloud banks in the summer, Lima then becomes unbearably hot as well as humid, and the sunshine seems to emphasize even more clearly the grimy buildings and lack of greenery in the central city.

THE LANDSCAPE

The city site. Lima sprawls well beyond its original Spanish site at a bridgeable point on the Río Rímac. Disgorging precipitously from the high Andes, the Rímac has formed a flat-topped alluvial cone, on which the early Spanish colonists established their settlement. Since

almost the entire coastal plain in central Peru consists of unconsolidated fluvio-glacial deposits, cliff erosion and earthquakes are continual threats. In expanding from its original site, the city has incorporated within its fabric various hills and valleys that are also prone to earth tremors and flash floods. One of the most notable characteristics of Lima is the barren, unvegetated desert that surrounds it on all sides; the grayish-yellow sands support almost no plant or animal life, save where water has been artificially provided.

Climate. Though Lima is located at a tropical latitude, the cool offshore Peru (also called Humboldt) Current helps produce a year-round temperate climate. Average temperature ranges from 60° to 64° F (16° to 18° C) in the winter months of May to November and 70° to 80° F (21° to 27° C) in the summer months of December to April. The cooling of the coastal air mass produces thick cloud cover throughout the winter, and the *garúa* (dense sea mist) often rolls in to blanket areas of the city. Precipitation, which rarely exceeds two inches (50 millimetres) per annum, usually results from the condensation of the *garúa*. Lima is perhaps best described as cold and damp in winter and hot and humid in summer.

Because clouds tend to trap airborne pollutants, Limeños can often taste the air. A permanent problem resulting from the high humidity is oxidation, rust being a common sight. Many of the wealthier citizens established winter homes on the coast north or south of the city proper or in such localities as La Molina, a short distance to the east of Lima, where the climate is free of fog and cloud.

The city layout. Lima contains a series of townscapes well-defined by its long history. The core of old Lima, delineated by Spanish colonists in the 16th century and partly enclosed by defensive walls in the 17th, retains its checkerboard street pattern. Bounded on the north by the Rímac and on the east, south, and west by broad avenues, old Lima contains a few restored colonial buildings (Torre Tagle Palace, the cathedral, and the Archbishop's Palace) interspersed among buildings of the 19th and 20th cen-

© Nat Norman—Rapho/Photo Researchers, Inc



The Archbishop's Palace, with colonial-style wooden balconies, on the Plaza de Armas.

The *garúa*

turies, many of which were built upon the sites of former colonial residences that had collapsed during the major earthquakes that have struck the city. The old walls, however, were demolished in the mid-19th century. The two principal squares (Plaza de Armas and Plaza Bolívar) still provide the foci of architectural interest within central Lima, and the enclosed wooden balconies so typical of the colonial city have now become features to be preserved or restored. The Presidential Palace (built on the site of Pizarro's house) and many other buildings reflect the past popularity of the French Empire style. On the north side of the Rímac, the old colonial suburb of the same name conserves relics of its past in its curved, narrow streets, tightly packed with single-story houses, and its Alameda de los Descalzos (Boulevard of the Barefoot Monks).

The former residential zone of central Lima has undergone several radical modifications, especially since the 1930s. Most of the old spacious mansions have been subdivided so that they now accommodate as many as 50 families. These inner-city slums (variously called *tugurios*, *corralones*, and *callejones*) have been occupied by immigrants from the countryside striving to gain a foothold in the urban economy and society. Sanitary conditions in such zones are often very poor.

Other parts of old Lima have experienced demolition and reconstruction. Housing has given way to banks, insurance offices, law firms, and government offices. Though there have been repeated attempts to stimulate pride in El Cercado (the formerly walled enclosure), most Limeños regard it as a place to pass through rather than to preserve and enhance. One finds little evidence of gentrification in Lima; unlike other Latin-American capitals and even other cities within Peru, central Lima contains relatively few outstanding architectural features.

Growth of
the city

Lima did not expand much beyond the walls of the old city until railways and trams were constructed in the mid-19th century. For the next 75 years growth was steady, the axes of urban development from old Lima assuming distinctive characters: the area west to Callao became the industrial corridor; the sweeping bay frontage to the south from Barranco to Magdalena was transformed into the choice residential zone; and eastward, toward Vitarte, a mix of industrial and lower-class suburbs sprang up. As the pace of urban expansion increased in the 1930s, small communities formed in the open country between Lima and the coast. These gradually coalesced into such urban districts as La Victoria, Lince, San Isidro, and Breña. The numerous farms and small tracts of cultivated land between suburbs and barren, dry land also became urbanized as immigrants from the interior occupied these areas. In the 1950s Lima became noted for these *barriadas* (squatter camps of shanties), which as they became more permanently established were renamed *pueblos jóvenes* ("young towns"). These communities have come to contain one-third of the population of metropolitan Lima. The older *pueblo jóvenes*, such as Comas, are now difficult to distinguish from the "established" sections of the city, since the early constructions of cardboard, tin cans, and wicker matting have long since given way to bricks, cement blocks, and neat gardens.

Lima's contemporary townscapes provide such contrasts that it is easy to forget that the rich and the poor belong to the same society. Within a few blocks one can move from luxury to abject poverty. With downtown Lima often heavily congested with traffic, suburban locations were chosen for many new businesses, factories, and shopping centres. The classic corner stores run by Chinese and Japanese immigrants and their descendants are fighting a losing battle against the competition of large, hygienic supermarkets, which in turn are increasingly encircled by *ambulantes* (street vendors).

THE PEOPLE

Just as the physical fabric of Lima has been transformed since the 1930s, so too has its population. It is now difficult to identify what might be called a true Limeño, for in a very real sense Lima has become the most Peruvian of cities; everywhere one can hear different accents, reflecting the myriad origins of the *provincianos* who have

made the city a microcosm of the nation. Before the arrival of the highland migrants (commonly called *serranos* or, if demonstrating what are perceived to be Indian characteristics, *cholos*), it was relatively easy to mark the difference between the white elite and other socioracial mixtures. Race, ethnicity, and class, however, now present a complexity that defies easy classification. The greatest difference that persists, and perhaps even increases, is that which divides the rich and influential from the poor and powerless. One has only to compare the elegance of those who stroll through Kennedy Park in Miraflores on a Saturday night with the squalor of those who beg in central Lima to realize that, in growing, the city has not developed. For the great majority of people access to piped water, sewage systems, inexpensive food, and steady employment are still dreams for the future.

The vast majority of Limeños are Roman Catholics, which gives the city a traditional, conservative atmosphere; this is evidenced by the enormous crowds of people who gather for such annual religious processions as El Señor de los Milagros ("the Lord of Miracles"), Santa Rosa de Lima, and San Martín de Porres. Growing numbers of residents from the slums and poor suburbs, however, have begun to question the church's positions on social and political issues, encouraged by priests who advocate what has become known as liberation theology.

Religion

THE ECONOMY

Whatever indicator is used to measure economic performance, Lima maintains a dominant position within Peru, accounting for three-fifths of the country's industrial output and nearly all of the volume of its financial transactions. The size of Lima's population makes it the premier market for all domestic and imported goods; Limeños make some four-fifths of the nation's consumer purchases each year.

Industry and commerce. Industry in Lima is located primarily in the old Callao-Lima-Vitarte corridor, with more recent additions in zones fringing the Pan-American Highway north and south of the city. Industrial activity is diverse, ranging from shipbuilding, automobile manufacturing, and oil refining to food processing and the manufacture of cement, chemicals, pharmaceuticals, plastics, textiles and clothing, and furniture. Much of this capital-intensive, heavily unionized industrial base, however, operates well below capacity, in most part because of the dire economic situation of Peru.

There has thus been a gradual de-emphasis of the more established industries, and since about 1970 a new type of informal, artisan-based industrial structure has developed. These small-scale, labour-intensive enterprises, which often are family controlled, have been better able to meet the demands of consumers by having goods more readily available (in part by avoiding bureaucratic red tape) and by offering goods for lower prices (in part by avoiding the payment of taxes).

Many industries have located within metropolitan Lima because of its pool of skilled labour, personal access to government officials, and the benefits of well-established networks of marketing and services such as banking. Manufacturing has not provided an adequate solution to the demands of the vast numbers who seek employment. One result has been the rapid rise in service jobs, the majority of which are informal in character. This type of employment has been estimated to account for at least 40 percent of total economic activity in the metropolitan area. The thousands of *ambulantes* who work the streets of central Lima have become a visual reminder of the lack of steady employment in the formal sector. One of the largest employers in Lima—directly and indirectly—is the national government. Its ministries, institutes, and other agencies provide jobs not only for an extensive bureaucracy but also for the hundreds of thousands of people who in various ways serve the needs of those fully employed.

Transportation. The railway line from Callao to Lima is the oldest in South America, while the line that climbs east past Vitarte and into the Andes reaches the highest point of any standard-gauge railway in the world. Other short lines radiate north to Paramonga and south to Luán.

The growth of automobile transportation has given rise to the heavily congested traffic conditions that exist in contemporary Lima. Although there is now a well-developed highway system in the metropolitan area, including an expressway between central Lima and Miraflores, the vast majority of Limeños must cope with an outdated street network and rely on three basic modes of transport: *colectivos*, shared taxis that can hold up to 10 passengers; small buses that can carry up to 20 persons; and municipal buses, of which more than half normally are broken down and the other half operate in bad repair.

Because transport in Lima is at best highly inefficient and at worst chaotic, hundreds of amateur taxi drivers (*piratas*), unlicensed and often ignorant of all but the most obvious locations within the city, offer their services to the harried or unwary pedestrian at peak traffic hours. Several plans for a subway system have been proposed for Lima, in part to overcome the obvious problems of the heavily congested and polluted centre but also to interconnect the peripheral suburbs more effectively and thus divert much traffic from the central city.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. The problems of controlling Lima's growth have proved difficult, but those of municipal administration have become almost insoluble. Metropolitan Lima consists of the provinces of Lima and Callao, the latter having departmental status, which are further divided into a total of 45 political districts. Each province and each district is administratively autonomous, so that citywide planning and development can be undertaken only by means of negotiated decisions. The capital district of Lima, with its long-established expertise in urban administration, has repeatedly called for the creation of a metropolitan authority that could more efficiently confront the many issues facing the region. Local district autonomy, however, which was won only after great political effort, has become a major obstacle to any unified approach, although a municipal law enacted in 1984 created a Metropolitan Council for Greater Lima (an assembly of district mayors) as well as agencies for improving cooperation between district councils and sharing technical assistance.

The system of generating and spending revenues in metropolitan Lima provides an example of the problems of interdistrict coordination. Since 1983 each district has been able both to generate its own revenues and to utilize them as it sees fit. Thus, there has been a growing disparity in the quality of services between the wealthy districts, which can generate adequate revenues for their needs, and the poor districts, which not only generate inadequate revenues but also are in most need of such services as water, sewers, electricity, and paved streets.

The differences in income and expenditures between rich and poor districts are, to some extent, paralleled by distinctive party affiliations and voting behaviour. The poorer districts have generally supported candidates from left-wing parties, while the more affluent suburbs have supported centre-right candidates. This interparty rivalry has hampered efforts at improving cooperation between districts as well as between the municipal and national government.

Services. The rapidity and scale of Lima's growth have placed great strains upon the provision of public services. Potable water, which in the past was obtained from the Rímac and from shallow local wells, now must be brought in via lakes and diverted rivers in the high Andes. Even so, it has been predicted that rising demand in Lima will make it necessary to find new water sources by the close of the 20th century. Equally difficult has been the provision of electricity. Only with the completion in the early 1970s of the expensive hydroelectric project on the Río Mantaro in Huancavelica department has affordable power been available for Lima's industry and residential population. These new sources of water and power, however, have been at the expense of the impoverished Andean departments that have provided them; the demand for payment for their resources by these departments has threatened to further strain Lima's municipal finances.

Within the capital itself the problems of providing ser-

VICES have been legion. Most municipalities have had barely enough income to finance their routine operations, with nothing left over to finance new projects. In addition, municipalities that have been able to allocate money for improved services often have been unable to adequately plan and execute what usually have been complex and highly technical projects. Finally, even when these projects have been built it has seldom been possible, given the penurious state of the majority of the population, to require payment for the actual cost of the services.

Caught between the need for inner-city renewal and suburban expansion, most municipalities have turned to the national government and such international agencies as the World Bank for assistance. Their argument has been that Lima's problems have become national problems and, as such, require national solutions. The response of the national government has been to point out that if the problems of the rural poor are not solved, especially in the departments of the southern Andes, more and more of those people will move to Lima, further compounding the city's difficulties.

CULTURAL LIFE

In spite of the many and complex problems that confront those who live in Lima, it is still the dominant and most vibrant cultural centre of Peru. Those who want to study almost any subject at the advanced level must go there. Lima contains the oldest and most distinguished universities in the country—including Universidad Nacional Mayor de San Marcos de Lima (National University of San Marcos), founded in 1551, and the Pontificia Universidad Católica del Perú (Pontific Catholic University of Peru), founded in 1917—as well as numerous other schools. Nearly all of the major academies, learned societies, and research institutes are located in metropolitan Lima, as are the national cultural institutions.

The numerous museums in the metropolitan area display the richness of Peru's pre-Columbian and colonial past. Within Lima itself are the well-restored burial sites (*huacas*) of the pre-Inca coastal cultures, and south of the city stand the remains of Pachacamac, one of Peru's largest pre-Hispanic cities. Dozens of other historic sites await funds for excavation and investigation, but almost all are threatened by urban construction.

Lima has several daily newspapers—*El Comercio*, founded in 1839, is the country's oldest—and numerous weekly periodicals, among which the magazine *Caretas* has become established as the newsweekly of Peru. Several television stations broadcast in colour, and some two dozen radio stations provide Lima with excellent coverage. Lima is not a city of bookstores or of book readers: the electronic media and a continual shortage of paper have combined to limit the circulation of the printed word. For the majority of lower-class Limeños the most popular reading materials are the comic books and dime novels that can be rented from street-corner stalls.

Recreation in Lima takes many forms, but perhaps no sports are more important than association football (soccer) for men and volleyball for women. Local soccer clubs have large and devoted followings. Other popular sports include horse racing, cockfighting, bullfighting, swimming, surfing, golf, tennis, and polo. Dozens of cinemas, theatre clubs, and discotheques provide nightlife, and there are scores of *peñas*, nightclubs featuring folk music. The music of Lima, as symbolized in the works of Chabuca Granda and Alicia Maguiño, is always popular, and it has enjoyed a renewed interest on the part of the public at large.

Both in the fashionable international-quality restaurants of central Lima and the bay area and in the hundreds of lesser cafés, *cevicherías*, and *picanterías* can be found a delicious variety of food. Fortunately for Lima, the migrants from other areas of Peru carried with them their highly flavoured regional dishes, making the city a gastronome's delight. Added to these foods are the excellent local beers, grape brandy (*pisco*), and inexpensive wines that are available.

One of the consequences of the massive migration to Lima has been the reinforcement of cultural ties between the capital's new urban communities and their localities

Structure

Media

of origin. Provincial and district clubs and associations celebrate weekly with songs, dances, and foods typical of the distinctive regions. Much of Peru's folklore can be learned in the heart of Lima itself. Yet the elite and the *cholo* alike share the pleasures of the beach and the ever-popular social equalizers, beer and *ceviche*, or *cebiche* (a typical coastal dish of marinated fish).

History

PRE-COLUMBIAN AND COLONIAL PERIODS

The area around Lima has been inhabited for thousands of years. Urban communities of significant size date from the pre-Inca Early Intermediate Period (c. 200 BC–AD 600), the most important being Pachacamac, which was an important religious site in both pre-Inca and Inca times. Much of the ransom demanded by the conquistador Francisco Pizarro for the Inca chief Atahualpa (Atahualpa) was obtained from Pachacamac.

Founding of Lima

The Spanish city of Lima was founded by Pizarro on Jan. 6, 1535, which, being Epiphany, prompted the name *Ciudad de los Reyes* ("City of Kings"). Although the name never stuck, Lima soon became the capital of the new viceroyalty of Peru, chosen over the old Inca capital of Cuzco to the southeast because the coastal location facilitated communication with Spain.

Lima developed into the centre of wealth and power for the entire viceroyalty; as the seat of the *audiencia* (high court), it administered royal justice; and, being the headquarters in the viceroyalty of the Inquisition, it pronounced on religious and moral matters. It also became the site of Peru's most prestigious associations and centres of learning, including the University of San Marcos (1551), the Peruvian Academy of Letters (1887), the National University of Engineering (1896), and the Pontific Catholic University of Peru (1917). José Hipólito Unzué founded a medical school there in 1808. From the late 17th to the mid-19th century, however, Lima grew extremely slowly in both area and population. The city was devastated by a powerful earthquake in 1746. Although it was rebuilt in grandiose fashion, influenced heavily by the European Enlightenment, it remained politically conservative and socially stratified. Lima maintained its loyalty during the struggles for Latin American independence in the early 19th century, with Peru becoming the last mainland colony to declare its independence from Spain (July 1821).

(D.J.Ro.)

THE MODERN CITY

Lima's development into a modern city began after the completion of the Lima-Callao railroad in 1851. Interurban railway links to Miraflores, Ancón, and Chosica followed in the next 20 years and provided the opportunity for suburban growth. The small, compact, pedestrian city gradually lost its wealthier residents, who physically distanced themselves from the lower classes by building mansions in and around Miraflores. Also during that period, Lima and Callao benefited from a boom in exports of nitrate-rich guano deposits, which were collected from islands off the Peruvian coast and shipped to Europe. However, Lima's prosperity subsequently declined as political turmoil swept the nation, and, as a result of the disastrous War of the Pacific, the Chilean military looted and occupied the city (1881–83), burning the National Library in the process.

Despite the loss of the library, the city's literary scene experienced a rebirth with Ricardo Palma's series of colonial legends and stories called *Tradiciones peruanas* ("Peruvian Traditions"), which appeared between 1872 and 1910. Influential literary figures of the early 20th century included the leftist political leader and essayist José Carlos Mariátegui and the poets César Vallejo, José María Eguren, and José Santos Chocano; although many of their works

focused on events outside of Lima (e.g., the plight of rural Indians), they exerted a profound influence on the intelligentsia of the city—and, by extension, the nation.

A new wave of urban expansion in the 1920s and '30s was set off by the automobile and the subsequent road-building program that improved transportation not only within the capital but also between Lima and other parts of the country. For the first time, migrants could reach Lima relatively easily, and this rich, powerful, and modernizing centre became a national magnet. The consequences for Lima were drastic. From 1940 to 1980 some two million people moved to the city. Hundreds of thousands of shanties were constructed on the bare, unoccupied slopes that rose above the red-tiled roofs of the inner suburbs and on the flat desert benches that encircled Lima. Individual acts of occupying unused and unclaimed pieces of land gave way to well-planned "invasions" involving many hundreds of new city residents. So enormous became the numbers of the self-help housing units that the government finally yielded to the residents' initiatives, awarding titles to the land and trying to provide basic services. Roughly one-third of metropolitan residents lived in *pueblos jóvenes* by 1990. A system of multilane expressways was built in the late 20th century to serve the city's expanding population, which had surpassed seven million by the late 1990s.

(D.J.Ro./Ed.)

Lima continues to influence nearly every facet of Peruvian national life—economic, political, and cultural. Since the mid-20th century, some of the more renowned works of novelist Mario Vargas Llosa have been set in Lima, including *La ciudad y los perros* (1963; "The City and the Dogs"; Eng. trans. *The Time of the Hero*) and *La tía Julia y el escribidor* (1977; *Aunt Julia and the Scriptwriter*). Among the more recent works focusing on Lima are Julio Ramón Ribeyro's tragicomic stories and Jaime Bayly's *Yo amo a mi mami* (1999; "I Love My Mom"), relating the experiences of a suburban child raised by household servants.

The city's historic centre was designated a UNESCO World Heritage site in 1988; in 1991 the site was redefined to include the former convent of San Francisco. However, Lima's historic buildings are threatened by elevated levels of air pollution from automobiles and buses, earthquakes (notably those of 1941 and 1970), and other hazards (such as a fire that destroyed the ornate municipal theatre in 1998). In the 1990s many of Lima's antique wooden balconies were repaired and restored.

BIBLIOGRAPHY. CAROLYN WALTON, *The City of Kings: A Guide to Lima* (1987), is a comprehensive guidebook. Basic information on Lima is also provided in *Lonely Planet: Peru*, 4th ed. (2000); and in *South American Handbook* (annual).

JUAN BROMLEY, *La fundación de la ciudad de los reyes* (1935); and JUAN BROMLEY and JOSÉ BARBAGELATA, *Evolución urbana de Lima . . .* (1945), provide the best introductions to the city's evolution up to the 20th century. Colonial Lima is described in MARK A. BURKHOLDER, *Politics of a Colonial Career: José Baquijano and the Audiencia of Lima*, 2nd ed. (1990); and JOSEPH DE MUGABURU and FRANCISCO DE MUGABURU, *Chronicle of Colonial Lima: The Diary of Josephe and Francisco Mugaburu, 1640–1697* (1975; originally published in Spanish, 1917–18). CHRISTINE HÜNEFELDT, *Paying the Price of Freedom: Family and Labor Among Lima's Slaves, 1800–1854* (1994), and *Liberalism in the Bedroom: Quarreling Spouses in Nineteenth-Century Lima* (2000), focus on gender, family relationships, and daily life.

Lima's growing pains since the mid-20th century are analyzed in SUSAN LOBO, *A House of My Own: Social Organization in the Squatter Settlements of Lima, Peru* (1982, reprinted 1992); PETER LLOYD, *The "Young Towns" of Lima: Aspects of Urbanization in Peru* (1980); HERNANDO DE SOTO, *The Other Path: The Invisible Revolution in the Third World* (1989; originally published in Spanish, 1986), which is a groundbreaking study of the informal economy (i.e., "black market") of Lima; and HENRY A. DIETZ, *Poverty and Problem-Solving Under Military Rule: The Urban Poor in Lima, Peru* (1980), and *Urban Poverty, Political Participation, and the State: Lima, 1970–1990* (1998). (Ed.)

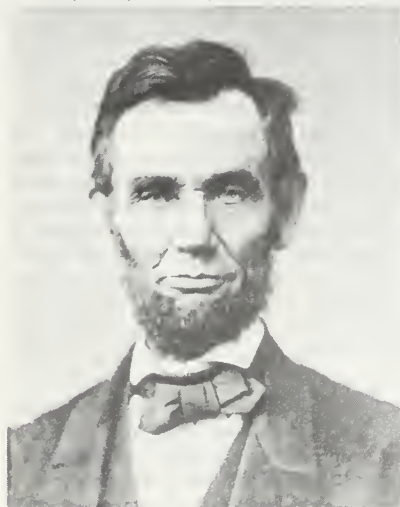
Influx of migrants and attendant problems

Lincoln

The 16th president of the United States (1861–65), Abraham Lincoln preserved the Union during the American Civil War and brought about the emancipation of the slaves.

Among American heroes, Lincoln continues to have a unique appeal for his fellow countrymen and also for people of other lands. This charm derives from his remarkable life story—the rise from humble origins, the dramatic death—and from his distinctively human and humane personality as well as from his historical role as saviour of the Union and emancipator of the slaves. His relevance endures and grows especially because of his eloquence as a spokesman for democracy. In his view, the Union was worth saving not only for its own sake but because it embodied an ideal, the ideal of self-government. In recent years, the political side to Lincoln's character, and his racial views in particular, have come under close scrutiny, as scholars continue to find him a rich subject for research.

By courtesy of the Library of Congress, Washington, D.C.



Lincoln, 1863.

CHILDHOOD AND YOUTH

Born on February 12, 1809, in a backwoods cabin three miles south of Hodgenville, Kentucky, Lincoln was two years old when he was taken to a farm in the neighbouring valley of Knob Creek. His father, Thomas Lincoln, was the descendant of an English weaver's apprentice who had migrated to Massachusetts in 1637. On June 12, 1806, he married Nancy Hanks, who appears to have been of illegitimate birth. She has been described as "stoop-shouldered, thin-breasted, sad," and fervently religious. Thomas and Nancy Lincoln had three children: Sarah, Abraham, and Thomas, who died in infancy.

In December 1816, faced with a lawsuit challenging the title to his Kentucky farm, Thomas Lincoln moved with his family to southwestern Indiana. There, as a squatter on public land, he hastily put up a "half-faced camp"—a crude structure of logs and boughs with one side open to the weather—in which the family took shelter behind a blazing fire. Soon he built a permanent cabin, and later he bought the land on which it stood. Abraham helped to clear the fields and to take care of the crops but early acquired a dislike for hunting and fishing. The unhappiest period of his boyhood followed the death of his mother in the autumn of 1818. As a ragged nine-year-old, he saw her buried in the forest, then faced a winter without the warmth of a mother's love. Fortunately, before the onset of a second winter, Thomas Lincoln brought home from Kentucky a new wife for himself and a new mother for the children. Sarah Bush Johnston Lincoln, a widow with two

Mother's death

girls and a boy of her own, had energy and affection to spare. She ran the household with an even hand, treating both sets of children as if she had borne them all; but she became especially fond of Abraham, and he of her. He afterward referred to her as his "angel mother."

His stepmother doubtless encouraged Lincoln's taste for reading, yet the original source of his desire to learn remains something of a mystery. Both his parents were almost completely illiterate, and he himself received little formal education. He once said that, as a boy, he had gone to school "by littles"—a little now and a little then—and his entire schooling amounted to no more than one year's attendance. According to his own statement, his early surroundings provided "absolutely nothing to excite ambition for education. Of course, when I came of age I did not know much. Still, somehow, I could read, write, and cipher to the rule of three; but that was all." Apparently the young Lincoln did not read a large number of books but thoroughly absorbed the few that he did read. These included Parson Weems's *Life and Memorable Actions of George Washington* (with its story of the little hatchet and the cherry tree), Daniel Defoe's *Robinson Crusoe*, John Bunyan's *Pilgrim's Progress*, and Aesop's *Fables*. From his earliest days he must have had some familiarity with the Bible, for it doubtless was the only book his family owned.

In March 1830 the Lincoln family undertook a second migration, this time to Illinois, with Lincoln himself driving the team of oxen. Having just reached the age of 21, he was six feet four inches tall, rawboned and lanky but muscular and physically powerful. He was especially noted for the skill and strength with which he could wield an ax. He spoke with a backwoods twang and walked in the long-striding, flat-footed, cautious manner of a plowman. Good-natured though somewhat moody, talented as a mimic and storyteller, he readily attracted friends. But he was yet to demonstrate whatever other abilities he possessed.

After his arrival in Illinois Lincoln tried his hand at a variety of occupations, including rail-splitter and flatboatman. In 1825 he settled in New Salem, where he worked as a storekeeper, a postmaster, and a surveyor. With the coming of the Black Hawk War (1832), he enlisted as a volunteer and was elected captain of his company. Meanwhile, aspiring to be a legislator, he was defeated in his first try and then repeatedly reelected to the state assembly. He considered blacksmithing as a trade but finally decided in favour of the law. Already having taught himself grammar and mathematics, he began to study law books. In 1836, having passed the bar examination, he began a law practice.

PRAIRIE LAWYER

The next year he moved to Springfield, Illinois, the new state capital, which offered many more opportunities for a lawyer than New Salem did. Eventually he became a partner of William H. Herndon, who was nearly 10 years younger than Lincoln, more widely read, more emotional at the bar, and generally more extreme in his views. Yet their partnership seems to have been as nearly perfect as such human arrangements ever are. Lincoln and Herndon kept few records of their law business and split the cash between them whenever either of them was paid. It seems they had no money quarrels.

Within a few years of his relocation to Springfield, Lincoln was earning from \$1,200 to \$1,500 annually (at that time the governor of the state received a salary of only \$1,200). To keep himself busy, he found it necessary not only to practice in the capital but also to follow the court as it made twice yearly rounds of its circuit through hundreds of miles of thinly settled prairie.

The coming of the railroads, especially after 1850, made

Migration to Illinois

Early legal career

travel easier and practice more remunerative. Lincoln served as a lobbyist for the Illinois Central Railroad, assisting it in getting a charter from the state, and thereafter he was a regular attorney for the company. He also handled cases for other railroads and for banks, insurance companies, and mercantile and manufacturing firms. His business included a number of patent suits and criminal trials. One of his most effective and famous pleas had to do with a murder case. A witness claimed that, by the light of the moon, he had seen Duff Armstrong, an acquaintance of Lincoln's, take part in a killing. Referring to an almanac for proof, Lincoln argued that the night had been too dark for the witness to have seen anything clearly, and with a sincere and moving appeal he won an acquittal.

By the time he began to be prominent in national politics, about 20 years after launching his legal career, Lincoln had made himself one of the most distinguished and successful lawyers in Illinois. He was noted not only for his shrewdness and practical common sense, which enabled him always to see to the heart of any legal case, but also for his invariable fairness and utter honesty.

PRIVATE LIFE

While residing in New Salem, Lincoln became acquainted with Ann Rutledge. Apparently he was fond of her, and certainly he grieved with the entire community at her untimely death, in 1835, at the age of 22. Afterward, stories were told of a grand romance between Lincoln and Rutledge, but these stories are not supported by sound historical evidence. A year after Rutledge's death, Lincoln carried on a halfhearted courtship with Mary Owens, who eventually concluded that he was "deficient in those little links which make up the chain of woman's happiness." She turned down his proposal.

So far as can be known, the first and only real love of Lincoln's life was Mary Todd. High-spirited, quick-witted, and well-educated, Todd came from a rather distinguished Kentucky family, and her Springfield relatives belonged to the social aristocracy of the town. Some of them frowned upon her association with Lincoln, and from time to time he too doubted that he could ever make her happy. Nevertheless, they became engaged. Then, on a day in 1841 that Lincoln recalled as the "fatal first of January," the engagement was broken, apparently on his initiative. For some time after that, Lincoln was overwhelmed by terrible depression and despondency. Finally the two were reconciled, and on November 4, 1842, they married.

Four children, all boys, were born to the Lincolns. Only the eldest, Robert Todd, survived to adulthood, though Lincoln's favourite, Thomas ("Tad"), who had a cleft palate and a lisp, outlived his father. Lincoln left the upbringing of his children largely to their mother, who was alternately strict and lenient in her treatment of them.

The Lincolns had a mutual affectionate interest in the doings and welfare of their boys, were fond of one another's company, and missed each other when apart, as existing letters show. Like most married couples, they also had domestic quarrels, which undoubtedly were exaggerated by contemporary gossips. She suffered from recurring headaches, fits of temper, and a sense of insecurity and loneliness that was intensified by her husband's long absences on the lawyer's circuit. After his election to the presidency, she was afflicted by the death of her son William ("Willie"), by the ironies of a war that made enemies of Kentucky relatives and friends, and by the unfair public criticisms of her as mistress of the White House. She developed an obsessive need to spend money, and she ran up embarrassing bills. She also staged some painful scenes of wifely jealousy. At last, in 1875, she was officially declared insane, though by that time she had undergone the shock of seeing her husband murdered at her side. During their earlier married life, she unquestionably encouraged her husband and served as a prod to his own ambition. During their later years together, she probably strengthened and tested his innate qualities of tolerance and patience.

With his wife, Lincoln attended Presbyterian services in Springfield and in Washington but never joined any church. Early in life Lincoln had been something of a skept-

tic and freethinker, and his reputation had been such that, as he once complained, the "church influence" was used against him in politics. When running for Congress in 1846, he issued a handbill to deny that he ever had "spoken with intentional disrespect of religion." He went on to explain that he had believed in the doctrine of necessity—"that is, that the human mind is impelled to action, or held in rest by some power over which the mind itself has no control." Throughout his life he also believed in dreams and other enigmatic signs and portents. As he grew older, and especially after he became president and faced the soul-troubling responsibilities of the Civil War, he developed a profound religious sense, and he increasingly personified necessity as God. He came to look upon himself quite humbly as an "instrument of Providence" and to view all history as God's enterprise.

Lincoln was fond of the Bible and knew it well. He also was fond of Shakespeare, and in private conversation he used many Shakespearean allusions, discussed problems of dramatic interpretation with considerable insight, and recited long passages from memory with rare feeling and understanding.

EARLY POLITICS

When Lincoln first entered politics, Andrew Jackson was president. Lincoln shared the sympathies that the Jacksonians professed for the common man, but he disagreed with the Jacksonian view that the government should be divorced from economic enterprise. "The legitimate object of government," he was later to say, "is to do for a community of people whatever they need to have done, but cannot do at all, or cannot do so well, for themselves, in their separate and individual capacities." Among the prominent politicians of his time, he most admired Henry Clay and Daniel Webster. Clay and Webster advocated using the powers of the federal government to encourage business and develop the country's resources by means of a national bank, a protective tariff, and a program of internal improvements for facilitating transportation. In Lincoln's view, Illinois and the West as a whole desperately needed such aid for economic development. From the outset, he associated himself with the party of Clay and Webster, the Whigs.

As a Whig member of the Illinois State Legislature, to which he was elected four times from 1834 to 1840, Lincoln devoted himself to a grandiose project for constructing with state funds a network of railroads, highways, and canals. Whigs and Democrats joined in passing an omnibus bill for these undertakings, but the panic of 1837 and the ensuing business depression brought about the abandonment of most of them.

While in the legislature he demonstrated that, though opposed to slavery, he was no abolitionist. In 1837, in response to the mob murder of Elijah Lovejoy, an antislavery newspaperman of Alton, the legislature introduced resolutions condemning abolitionist societies and defending slavery within the Southern states as "sacred" by virtue of the federal Constitution. Lincoln refused to vote for the resolutions. Together with a fellow member, he drew up a protest that declared, on the one hand, that slavery was "founded on both injustice and bad policy" and, on the other, that "the promulgation of abolition doctrines tends rather to increase than to abate its evils."

During his single term in Congress (1847-49), Lincoln, as the lone Whig from Illinois, gave little attention to legislative matters. He proposed a bill for the gradual and compensated emancipation of slaves in the District of Columbia, but, because it was to take effect only with the approval of the "free white citizens" of the district, it displeased abolitionists as well as slaveholders and never was seriously considered.

Lincoln devoted much of his time to presidential politics—to unmaking one president, a Democrat, and making another, a Whig. He found an issue and a candidate in the Mexican War. With his "spot resolutions," he challenged the statement of President James K. Polk that Mexico had started the war by shedding American blood upon American soil. Along with other members of his party, Lincoln voted to condemn Polk and the war while also voting for

Marriage to Mary Todd

Illinois state legislator

Religious views

supplies to carry it on. At the same time, he laboured for the nomination and election of the war hero Zachary Taylor. After Taylor's success at the polls, Lincoln expected to be named commissioner of the general land office as a reward for his campaign services, and he was bitterly disappointed when he failed to get the job. His criticisms of the war, meanwhile, had not been popular among the voters in his own congressional district. At the age of 40, frustrated in politics, he seemed to be at the end of his public career.

THE ROAD TO THE PRESIDENCY

For about five years Lincoln took little part in politics, and then a new sectional crisis gave him a chance to reemerge and rise to statesmanship. In 1854 his political rival Stephen A. Douglas maneuvered through Congress a bill for reopening the entire Louisiana Purchase to slavery and allowing the settlers of Kansas and Nebraska (with "popular sovereignty") to decide for themselves whether to permit slaveholding in those territories. The Kansas-Nebraska Act provoked violent opposition in Illinois and the other states of the old Northwest. It gave rise to the Republican Party while speeding the Whig Party on its way to disintegration. Along with many thousands of other homeless Whigs, Lincoln soon became a Republican (1856).

Lincoln challenged the incumbent Douglas for the Senate seat in 1858, and the series of debates they engaged in throughout Illinois was political oratory of the highest order. Both men were shrewd debaters and accomplished stump speakers, though they could hardly have been more different in style and appearance—the short and pudgy Douglas, whose stentorian voice and graceful gestures swayed audiences, and the tall, homely, almost emaciated-looking Lincoln, who moved awkwardly and whose voice was piercing and shrill. Lincoln's prose and speeches, however, were eloquent, pithy, powerful, and free of the verbosity so common in his day. The debates were published in 1860, together with a biography of Lincoln, in a best-selling book that Lincoln himself compiled and marketed as part of his campaign.

In their basic views, Lincoln and Douglas were not as far apart as they seemed in the heat of political argument. Neither was abolitionist or proslavery. But Lincoln, unlike Douglas, insisted that Congress must exclude slavery from the territories. He disagreed with Douglas's belief that the territories were by nature unsuited to the slave economy and that no congressional legislation was needed to prevent the spread of slavery into them. In one of his most famous speeches he said: "A house divided against itself cannot stand. I believe the government cannot endure permanently half slave and half free." He predicted that the country eventually would become "all one thing, or all the other." He agreed with Thomas Jefferson and other founding fathers, however, that slavery should be merely contained, not directly attacked. In fact, when politically expedient, he reassured his audiences that he did not endorse citizenship for Negroes or believe in the equality of the races. "I am not, nor ever have been, in favour of bringing about in any way the social and political equality of the white and black races," he told a crowd in Charleston, Illinois. "I am not nor ever have been in favour of making voters or jurors of Negroes, nor of qualifying them to hold office, nor to intermarry with white people." There is, he added, "a physical difference between the white and black races which I believe will forever forbid the two races living together on terms of social and political equality."

In the end, Lincoln lost the election to Douglas. Although the outcome did not surprise him, it depressed him deeply. Lincoln had, nevertheless, gained national recognition and soon began to be mentioned as a presidential prospect for 1860.

On May 18, 1860, after Lincoln and his friends had made skillful preparations, he was nominated on the third ballot at the Republican National Convention in Chicago. He then put aside his law practice and, though making no stump speeches, gave full time to the direction of his campaign. With the Republicans united, the Democrats divided, and a total of four candidates in the field, he carried the election on November 6. Although he received no votes from the Deep South and only 40 out of 100 in the coun-

try as a whole, the popular votes were so distributed that he won a clear and decisive majority in the electoral college.

PRESIDENT LINCOLN

After Lincoln's election and before his inauguration, the state of South Carolina proclaimed its withdrawal from the Union. To forestall similar action by other Southern states, various compromises were proposed in Congress. The most important, the Crittenden Compromise, included constitutional amendments guaranteeing slavery forever in the states where it already existed and dividing the territories between slavery and freedom. Although Lincoln had no objection to the first of these amendments, he was unalterably opposed to the second, and indeed to any scheme infringing in the slightest upon the free-soil plank of his party's platform. He feared that a territorial division, by sanctioning the principle of slavery extension, would only encourage planter imperialists to seek new slave territory south of the American border and thus would "put us again on the highroad to a slave empire." From his home in Springfield he advised Republicans in Congress to vote against a division of the territories, and the proposal was killed in committee. Six additional states then seceded and, with South Carolina, combined to form the Confederate States of America.

Thus, before Lincoln had even moved into the White House, a disunion crisis was upon the country. Attention, North and South, focused upon Fort Sumter, in Charleston Harbor, South Carolina. This fort, still under construction, was garrisoned by U.S. troops under Major Robert Anderson. The Confederacy claimed it and, from other harbour fortifications, threatened it. Foreseeing trouble, Lincoln, while still in Springfield, confidentially requested Winfield Scott, general in chief of the U.S. Army, to be prepared "to either hold, or retake, the forts, as the case may require, at, and after the inauguration." In his inaugural address (March 4, 1861), besides upholding the Union's indestructibility and appealing for sectional harmony, Lincoln restated his Sumter policy as follows:

The power confided to me, will be used to hold, occupy, and possess the property, and places belonging to the government, and to collect the duties and imposts; but beyond what may be necessary for these objects, there will be no invasion—no using of force against, or among the people anywhere.

Then, near the end, addressing the absent Southerners: "You can have no conflict, without being yourselves the aggressors."

Outbreak of war. No sooner was he in office than Lincoln received word that the Sumter garrison, unless supplied or withdrawn, would shortly be starved out. Still, for about a month, Lincoln delayed acting. He was beset by contradictory advice. General Scott, Secretary of State William H. Seward, and others urged him to abandon the fort; and Seward, through a go-between, gave a group of Confederate commissioners to understand that the fort would in fact be abandoned. But many Republicans insisted that any show of weakness would bring disaster to their party and to the Union. Finally, Lincoln ordered the preparation of two relief expeditions, one for Fort Sumter and the other for Fort Pickens, in Florida. (He afterward said he would have been willing to withdraw from Sumter if he could have been sure of holding Pickens.) Before the Sumter expedition, he sent a messenger to tell the South Carolina governor:

I am directed by the President of the United States to notify you to expect an attempt will be made to supply Fort-Sumter [sic] with provisions only; and that, if such attempt be not resisted, no effort to throw in men, arms, or ammunition, will be made, without further notice, or in case of an attack upon the Fort.

Without waiting for the arrival of Lincoln's expedition, the Confederate authorities presented to Major Anderson a demand for Sumter's prompt evacuation, which he refused. On April 12, 1861, at dawn, the Confederate batteries in the harbour opened fire.

"Then, and thereby," Lincoln informed Congress when it met on July 4, "the assailants of the Government, began the conflict of arms." The Confederates, however, accused him of being the real aggressor. They said he had cleverly

Lincoln-Douglas debates

Secession of Confederate states

Election as president

Attack on Fort Sumter

maneuvered them into firing the first shot so as to put upon them the onus of war guilt. Although some historians have repeated this charge, it appears to be a gross distortion of the facts. Lincoln was determined to preserve the Union, and to do so he thought he must take a stand against the Confederacy. He concluded he might as well take this stand at Sumter.

Lincoln's primary aim was neither to provoke war nor to maintain peace. In preserving the Union, he would have been glad to preserve the peace also, but he was ready to risk a war that he thought would be short.

After the firing on Fort Sumter, Lincoln called upon the state governors for troops (Virginia and three other states of the upper South responded by joining the Confederacy). He then proclaimed a blockade of the Southern ports. These steps—the Sumter expedition, the call for volunteers, and the blockade—were the first important decisions of Lincoln as commander in chief of the army and navy. But he still needed a strategic plan and a command system for carrying it out.

General Scott advised him to avoid battle with the Confederate forces in Virginia, to get control of the Mississippi River, and by tightening the blockade to hold the South in a gigantic squeeze. Lincoln had little confidence in Scott's comparatively passive and bloodless "Anaconda" plan. He believed the war must be actively fought if it ever was to be won. Overruling Scott, he ordered a direct advance on the Virginia front, which resulted in defeat and rout for the federal forces at Bull Run (July 21, 1861). After a succession of more or less sleepless nights, Lincoln produced a set of memorandums on military policy. His basic thought was that the armies should advance concurrently on several fronts and should move so as to hold and use the support of Unionists in Missouri, Kentucky, western Virginia, and eastern Tennessee. As he later explained:

I state my general idea of this war to be that we have the greater numbers, and the enemy has the greater facility of concentrating forces upon points of collision; that we must fail, unless we can find some way of making our advantage an over-match for his; and that this can only be done by menacing him with superior forces at different points, at the same time.

This, with the naval blockade, comprised the essence of Lincoln's strategy.

Leadership in war. As a war leader, Lincoln employed the style that had served him as a politician—a description of himself, incidentally, that he was not ashamed to accept. He preferred to react to problems and to the circumstances that others had created rather than to originate policies and lay out long-range designs. In candour he would write: "I claim not to have controlled events, but confess plainly that events have controlled me." His guiding rule was: "My policy is to have no policy." He was not unprincipled, but he was a practical man, mentally nimble and flexible, and if one action or decision proved unsatisfactory in practice he was willing to experiment with another.

From 1861 to 1864, while hesitating to impose his ideas upon his generals, Lincoln experimented with command personnel and organization. Accepting the resignation of Scott (November 1861), he put George B. McClellan in charge of the armies as a whole. After a few months, disgusted by the slowness of McClellan, he demoted him to the command of the Army of the Potomac alone. He then tried a succession of commanders for the army in Virginia but was disappointed with each of them in turn. Meanwhile, he had in Henry W. Halleck a general in chief who shrank from making important decisions. For nearly two years the Federal armies lacked effective unity of command. Lincoln, besides transmitting official orders through Halleck, also communicated directly with the generals, sending personal suggestions in his own name. To generals opposing Robert E. Lee, he suggested that the object was to destroy Lee's army, not to capture Richmond or to drive the invader from Northern soil.

Finally Lincoln looked to the West for a top general. He admired the Vicksburg Campaign of Ulysses S. Grant in Mississippi, and, nine days after the Vicksburg surrender (which occurred on July 4, 1863), he sent Grant a "grateful acknowledgment for the almost inestimable service" he had done the country. In March 1864 Lincoln promoted

Grant to lieutenant general and gave him command of all the federal armies. At last Lincoln had found a man who, with such able subordinates as William T. Sherman, Philip Sheridan, and George H. Thomas, could put into effect those parts of Lincoln's concept of a large-scale, coordinated offensive that still remained to be carried out. Grant was only a member, though an important one, of a top-command arrangement that Lincoln eventually had devised. Overseeing everything was Lincoln himself, the commander in chief. Taking the responsibility for men and supplies was Secretary of War Edwin M. Stanton. Serving as a presidential adviser and as a liaison with military men was Halleck, the chief of staff. And directing all the armies, while accompanying Meade's Army of the Potomac, was Grant, the general in chief.

Lincoln combined statecraft and the overall direction of armies with an effectiveness that increased year by year. His achievement is all the more remarkable in view of his lack of training and experience in the art of warfare. This lack may have been an advantage as well as a handicap. Unhindered by outworn military dogma, Lincoln could all the better apply his practical insight and common sense—some would say his military genius—to the winning of the Civil War.

There can be no doubt of Lincoln's deep and sincere devotion to the cause of personal freedom. Before his election to the presidency he had spoken often and eloquently on the subject. In 1854, for example, he said he hated the Douglas attitude of indifference toward the possible spread of slavery to new areas. "I hate it because of the monstrous injustice of slavery itself," he declared. "I hate it because it deprives our republican example of its just influence in the world; enables the enemies of free institutions with plausibility to taunt us as hypocrites."

Yet, as president, Lincoln was at first reluctant to adopt an abolitionist policy. There were several reasons for his hesitancy. He had been elected on a platform pledging no interference with slavery within the states, and in any case he doubted the constitutionality of federal action under the circumstances. He was concerned about the possible difficulties of incorporating nearly four million Negroes, once they had been freed, into the nation's social and political life. Above all, he felt that he must hold the border slave states in the Union, and he feared that an abolitionist program might impel them, in particular his native Kentucky, toward the Confederacy. So he held back while others went ahead. When General John C. Frémont and General David Hunter, within their respective military departments, proclaimed freedom for the slaves of disloyal masters, Lincoln revoked the proclamations. When Congress passed confiscation acts in 1861 and 1862, he refrained from a full enforcement of the provisions authorizing him to seize slave property. And when Horace Greeley in the *New York Tribune* appealed to him to enforce these laws, Lincoln patiently replied (August 22, 1862):

My paramount object in this struggle is to save the Union, and is not either to save or to destroy slavery. If I could save the Union without freeing any slave I would do it; and if I could save it by freeing all the slaves I would do it; and if I could save it by freeing some and leaving others alone, I would also do that.

Meanwhile, in response to rising antislavery sentiment, Lincoln came forth with an emancipation plan of his own. According to his proposal, the slaves were to be freed by state action, the slaveholders were to be compensated, the federal government was to share the financial burden, the emancipation process was to be gradual, and the freedmen were to be colonized abroad. Congress indicated its willingness to vote the necessary funds for the Lincoln plan, but none of the border slave states were willing to launch it, and in any case few Negro leaders desired to see their people sent abroad.

While still hoping for the eventual success of his gradual plan, Lincoln took quite a different step by issuing his preliminary (September 22, 1862) and his final (January 1, 1863) Emancipation Proclamation. This famous decree, which he justified as an exercise of the president's war powers, applied only to those parts of the country actually under Confederate control, not to the loyal slave states nor to the federally occupied areas of the Confederacy. Direct-

Lincoln
and
slavery

Strategy
as com-
mander

Emanci-
pation
Procla-
mation

The
Thirteenth
Amend-
ment

ly or indirectly, the proclamation brought freedom during the war to fewer than 200,000 slaves. Yet it had great significance as a symbol. It indicated that the Lincoln government had added freedom to reunion as a war aim, and it attracted liberal opinion in England and Europe to increased support of the Union cause.

Lincoln himself doubted the constitutionality of his step, except as a temporary war measure. After the war, the slaves freed by the proclamation would have risked re-enslavement had nothing else been done to confirm their liberty. But something else was done: the Thirteenth Amendment was added to the Constitution, and Lincoln played a large part in bringing about this change in the fundamental law. Through the chairman of the Republican National Committee he urged the party to include a plank for such an amendment in its platform of 1864. The plank, as adopted, stated that slavery was the cause of the rebellion, that the president's proclamation had aimed "a death blow at this gigantic evil," and that a constitutional amendment was necessary to "terminate and forever prohibit" it. When Lincoln was reelected on this platform and the Republican majority in Congress was increased, he was justified in feeling, as he apparently did, that he had a mandate from the people for the Thirteenth Amendment. The newly chosen Congress, with its overwhelming Republican majority, was not to meet until after the lame duck session of the old Congress during the winter of 1864-65. Lincoln did not wait. Using his resources of patronage and persuasion upon certain of the Democrats, he managed to get the necessary two-thirds vote before the session's end. He rejoiced as the amendment went out to the states for ratification, and he rejoiced again and again as his own Illinois led off and other states followed one by one in acting favourably upon it. (He did not live to rejoice in its ultimate adoption.)

Lincoln deserves his reputation as the Great Emancipator. His claim to that honour, if it rests uncertainly upon his famous proclamation, has a sound basis in the support he gave to the antislavery amendment. It is well founded also in his greatness as the war leader who carried the nation safely through the four-year struggle that brought freedom in its train. And, finally, it is strengthened by the practical demonstrations he gave of respect for human worth and dignity, regardless of colour. During the last two years of his life he welcomed Negroes as visitors and friends in a way no president had done before. One of his friends was the distinguished former slave Frederick Douglass, who wrote: "In all my interviews with Mr. Lincoln I was impressed with his entire freedom from prejudice against the colored race."

Wartime politics. To win the war, President Lincoln had to have popular support. The reunion of North and South required, first of all, a certain degree of unity in the North. But the North contained various groups with special interests of their own. Lincoln faced the task of attracting to his administration the support of as many divergent groups and individuals as possible. Fortunately for the Union cause, he was a president with rare political skill. He had the knack of appealing to fellow politicians and talking to them in their own language. He had a talent for smoothing over personal differences and holding the loyalty of men antagonistic to one another. Inheriting the spoils system, he made good use of it, disposing of government jobs so as to strengthen his administration and further its official aims.

The opposition party remained alive and strong. Its membership included war Democrats and peace Democrats, often called "Copperheads," a few of whom collaborated with the enemy. Lincoln did what he could to cultivate the assistance of the war Democrats, as in securing from Congress the timely approval of the Thirteenth Amendment. So far as feasible, he conciliated the peace Democrats. In dealing with persons suspected of treasonable intent, Lincoln at times authorized his generals to make arbitrary arrests. He justified this action on the ground that he had to allow some temporary sacrifice of parts of the Constitution in order to maintain the Union and thus preserve the Constitution as a whole. He let his generals suspend several newspapers, but only for short periods, and he promptly

revoked a military order suppressing the hostile *Chicago Times*.

Considering the dangers and provocations of the time, Lincoln was quite liberal in his treatment of political opponents and the opposition press. He was by no means the dictator critics often accused him of being. Nevertheless, his abrogating of civil liberties, especially his suspension of the privilege of the writ of habeas corpus, disturbed Democrats, Republicans, and even members of his own cabinet. In the opinion of a soldier from Massachusetts, the president, "without the people having any legal means to prevent it, is only prevented from exercising a Russian despotism by the fear he may have of shocking too much the sense of decency of the whole world." Even Lincoln's friend Orville Hickman Browning believed the arrests ordered by the president were "illegal and arbitrary, and did more harm than good, weakening instead of strengthening the government." Yet Lincoln defended his actions, arguing that the Constitution provided for the suspension of such liberties "in cases of Rebellion or Invasion, [when] the public Safety may require it." Moreover, posed Lincoln with rhetorical flare, "Must I shoot a simpleminded soldier boy who deserts" and "not touch a hair of a wily agitator who induces him to desert?"

Within his own party, Lincoln confronted factional divisions and personal rivalries that caused him as much trouble as did the activities of the Democrats, though he and most of his fellow partisans agreed fairly well upon their principal economic aims. With his approval, the Republicans enacted into law the essentials of the program he had advocated from his early Whig days—a protective tariff; a national banking system; and federal aid for internal improvements, in particular for the construction of a railroad to the Pacific Coast. The Republicans disagreed among themselves, however, on many matters regarding the conduct and purposes of the war. Two main factions arose: the "Radicals" and the "Conservatives." Lincoln himself inclined in spirit toward the Conservatives, but he had friends among the Radicals as well, and he strove to maintain his leadership over both. In appointing his cabinet, he chose his several rivals for the 1860 nomination and, all together, gave representation to every important party group. Wisely he included the outstanding Conservative, Seward, and the outstanding Radical, Salmon P. Chase. Cleverly he overcame cabinet crises and kept these two opposites among his official advisers until Chase's resignation in 1864.

Lincoln had to deal with even more serious factional uprisings in Congress. The big issue was the "reconstruction" of the South. The seceded states of Louisiana, Arkansas, and Tennessee having been largely recovered by the federal armies, Lincoln late in 1863 proposed his "ten percent plan," according to which new state governments might be formed when 10 percent of the qualified voters had taken an oath of future loyalty to the United States. The Radicals rejected Lincoln's proposal as too lenient, and they carried through Congress the Wade-Davis Bill, which would have permitted the remaking and readmission of states only after a majority had taken the loyalty oath. When Lincoln pocket-vetoed that bill, its authors published a "manifesto" denouncing him.

Lincoln was already the candidate of the "Union" (that is, the Republican) party for reelection to the presidency, and the Wade-Davis manifesto signalized a movement within the party to displace him as the party's nominee. But even after the movement collapsed, the party remained badly divided. A rival Republican candidate, John C. Frémont, nominated much earlier by a splinter group, was still in the field. Leading Radicals promised to procure Frémont's withdrawal if Lincoln would obtain the resignation of his conservative postmaster general, Montgomery Blair. Eventually Frémont withdrew and Blair resigned. The party was reunited in time for the election of 1864.

In 1864, as in 1860, Lincoln was the chief strategist of his own electoral campaign. He took a hand in the management of the Republican Speakers' Bureau, advised state committees on campaign tactics, hired and fired government employees to strengthen party support, and did his best to enable as many soldiers and sailors as possible to

Abrogation
of civil
libertiesRecon-
struction

vote. Most of the citizens in uniform voted Republican. He was reelected with a large popular majority (55 percent) over his Democratic opponent, General George B. McClellan.

In 1864 the Democratic platform called for an armistice and a peace conference, and prominent Republicans as well as Democrats demanded that Lincoln heed Confederate peace offers, irregular and illusory though they were. In a public letter, he stated his own conditions:

Any proposition which embraces the restoration of peace, the integrity of the whole Union, and the abandonment of slavery, and which comes by and with an authority that can control the armies now at war against the United States will be received and considered by the Executive government of the United States, and will be met by liberal terms on other substantial and collateral points.

When Conservatives protested to him against the implication that the war must go on to free the slaves, even after reunion had been won, he explained, "To me it seems plain that saying reunion and abandonment of slavery would be considered, if offered, is not saying that nothing else or less would be considered, if offered." After his reelection, in his annual message to Congress, he said, "In stating a single condition of peace, I mean simply to say that the war will cease on the part of the government, whenever it shall have ceased on the part of those who began it." On February 3, 1865, he met personally with Confederate commissioners on a steamship in Hampton Roads, Virginia. He promised to be liberal with pardons if the South would quit the war, but he insisted on reunion as a precondition for any peace arrangement. In his second inaugural address he embodied the spirit of his policy in the famous words "with malice toward none; with charity for all." His terms satisfied neither the Confederate leaders nor the Radical Republicans, and so no peace was possible until the final defeat of the Confederacy.

Postwar policy. At the end of the war, Lincoln's policy for the defeated South was not clear in all its details, though he continued to believe that the main object should be to restore the "seceded States, so-called," to their "proper practical relation" with the Union as soon as possible. He possessed no fixed and uniform program for the region as a whole. With respect to states like Louisiana and Tennessee, he continued to urge acceptance of new governments set up under his 10 percent plan during the war. With respect to states like Virginia and North Carolina, he seemed willing to use the old rebel governments temporarily as a means of transition from war to peace. He was on record as opposing the appointment of "strangers" (carpetbaggers) to govern the South. A program of education for the freedmen, he thought, was essential to preparing them for their new status. He also suggested that the vote be given immediately to some Negroes—"as, for instance, the very intelligent, and especially those who have fought gallantly in our ranks."

On the question of reconstruction, however, Lincoln and the extremists of his own party stood even farther apart in early 1865 than a year before. Some of the Radicals were beginning to demand a period of military occupation for the South, the confiscation of planter estates and their division among the freedmen, and the transfer of political power from the planters to their former slaves. In April 1865 Lincoln began to modify his own stand in some respects and thus to narrow the gap between himself and the Radicals. He recalled the permission he had given for the assembling of the rebel legislature of Virginia, and he approved in principle—or at least did not disapprove—Stanton's scheme for the military occupation of Southern states. After the cabinet meeting of April 14, Attorney General James Speed inferred that Lincoln was moving toward the Radical position. What Lincoln's reconstruction policy would have been if he had lived to complete his second term can only be guessed at.

On the evening of April 14, 1865, 26-year-old John Wilkes Booth—a rabid advocate of slavery with ties to the South and the flamboyant son of one of the most distinguished theatrical families of the 19th century—shot Lincoln as he sat in Ford's Theatre in Washington. Early the next morning Lincoln died.

(R.N.C./Ed.)

REPUTATION AND CHARACTER

"Now he belongs to the ages," Stanton is supposed to have said as Lincoln took his last breath. Many thought of Lincoln as a martyr. The assassination had occurred on Good Friday, and on the following Sunday, memorable as "Black Easter," hundreds of speakers found a sermon in the event. Some of them saw more than mere chance in the fact that assassination day was also crucifixion day. One declared, "Jesus Christ died for the world; Abraham Lincoln died for his country." Thus the posthumous growth of his reputation was influenced by the timing and circumstances of his death, which won for him a kind of sainthood.

Among the many who remembered Lincoln from personal acquaintance, William Herndon, his former law partner, was sure he had known him more intimately than any of the rest and influenced the world's conception of him more than all the others put together. When Lincoln died, Herndon began a new career as Lincoln authority, collecting reminiscences wherever he could find them and adding his own store of memories. Although admiring Lincoln, he objected to the trend toward sanctifying him. He saw, as the main feature of Lincoln's life, the far more than ordinary rise of a self-made man, a rise from the lowest depths to the greatest heights. To emphasize this point, Herndon gave his most eager attention to evidence of the dismal and sordid in Lincoln's background. An extremely significant event in Lincoln's development, as Herndon viewed it, was a "romance of much reality" with Ann Rutledge. Lincoln loved no one but Rutledge and, after her death, never ceased to grieve for her. His memory of her both saddened and inspired him. As for Mary Todd, she married Lincoln out of spite, then devoted herself to making him miserable. So Herndon would have it, and after him countless biographers, novelists, and playwrights elaborated upon his views, which persist as accepted knowledge about Lincoln despite their refutation by historical scholarship.

Lincoln has become a myth as well as a man. The Lincoln of legend has grown into a protean god who can assume a shape to please almost anyone. He is Old Abe and at the same time a natural gentleman. He is Honest Abe and yet a being of superhuman shrewdness and cunning. He is also Father Abraham, the wielder of authority, the support of the weak; and he is an equal, a neighbour, and a friend. But there is a malevolent Lincoln as well, and to many Southerners from the time of the Civil War and to some conservative critics today, Lincoln is the wicked slayer of liberty and states' rights and the father of the all-controlling national state.

Lincoln's reputation began to grow while he was still alive. In the midst of the Civil War, for instance, the *Washington Chronicle* found a resemblance between him and George Washington in their "sure judgment," "perfect balance of thoroughly sound faculties," and "great calmness of temper, great firmness of purpose, supreme moral principle, and intense patriotism."

Lincoln's best ideas and finest phrases were considered and written and rewritten with meticulous revisions. Some resulted from a slow gestation of thought and phrase through many years. One of his themes—probably his central theme—was the promise and the problem of self-government. As early as 1838, speaking to the Young Men's Lyceum of Springfield on "The Perpetuation of Our Political Institutions," he recalled the devotion of his Revolutionary forefathers to the cause and went on to say:

Their ambition aspired to display before an admiring world, a practical demonstration of the truth of a proposition, which had hitherto been considered, at best no better, than problematical; namely, the capability of a people to govern themselves.

Again and again he returned to this idea, especially after the coming of the Civil War, and he steadily improved his phrasing. Finally at Gettysburg, Pennsylvania, he made the culminating, supreme statement, concluding with the words:

... that from these honored dead we take increased devotion to that cause for which they gave the last full measure of devotion—that we here highly resolve that these dead shall not have died in vain—that this nation, under God, shall have a new birth of freedom—and that government of the people, by the people, for the people, shall not perish from the earth.

The myth of Lincoln

BIBLIOGRAPHY. A guide for the general reader is PAUL M. ANGLE, *A Shelf of Lincoln Books: A Critical, Selective Bibliography of Lincolniana* (1946, reissued 1972). Practically all the known writings of Lincoln himself are available in *The Collected Works of Abraham Lincoln*, ed. by ROY P. BASLER, 9 vol. (1953–55), with two supplements (1974, 1990). A judicious selection from these volumes is reprinted in *The Living Lincoln*, ed. by PAUL M. ANGLE and EARL SCHENK MIERS (1955, reissued 1992). *Lincoln on Democracy*, ed. by MARIO M. CUOMO and HAROLD HOLZER (1990), contains Lincoln's writings on this subject. THOMAS S. SCHWARTZ (ed.), *For a Vast Future Also: Essays from the Journal of the Abraham Lincoln Association* (1999), includes essays on emancipation and Lincoln's legacy.

Classic multivolume biographies are JOHN G. NICOLAY and JOHN HAY, *Abraham Lincoln: A History*, 10 vol. (1890, reissued 1917), also available in an abridged ed. edited by PAUL M. ANGLE, 1 vol. (1966); ALBERT BEVERIDGE, *Abraham Lincoln, 1809–1858*, 2 vol. (1928, reissued 1971); CARL SANDBURG, *Abraham Lincoln: The Prairie Years*, 2 vol. (1926), and *Abraham Lincoln: The War Years*, 4 vol. (1939), both reissued together in 1 vol. (1984); and J.G. RANDALL, *Lincoln, the President*, 4 vol. (1945–55). Indispensable for the politics of the 1850s is ALLAN NEVINS, *The Emergence of Lincoln*, 2 vol. (1950). One-volume biographies include BENJAMIN P. THOMAS, *Abraham Lincoln* (1952, reissued 1986); STEPHEN B. OATES, *With Malice Toward None* (1977, reissued 1985), and *Abraham Lincoln, the Man Behind the Myths* (1984); OSCAR HANDLIN and LILIAN HANDLIN, *Abraham Lincoln and the Union* (1980); PHILIP B. KUNHARDT, JR., PHILIP B. KUNHARDT III, and PETER W. KUNHARDT, *Lincoln* (1992), containing 900 pictures; MICHAEL BURLINGAME, *The Inner World of Abraham Lincoln* (1994), a psychobiography; DAVID HERBERT DONALD, *Lincoln* (1995); ALLEN C. GUELZO, *Abraham Lincoln: Redeemer President* (1999); and WILLIAM LEE MILLER, *Lincoln's Virtues: An Ethical Biography* (2002). PAUL HORGAN, *Citizen of New Salem* (also published as *Abraham Lincoln, Citizen of New Salem*, 1961), concentrates on Lincoln's early life; while MARK E. NEELY, JR., *The Last Best Hope of Earth: Abraham Lincoln and the Promise of America* (1993), traces Lincoln's later political life through his own speeches. A collection of valuable appreciations is found in RALPH G. NEWMAN (ed.), *Lincoln for the Ages* (1960).

Matters of controversy may be found in LLOYD LEWIS, *Myths after Lincoln* (1929, reissued as *The Assassination of Lincoln: History and Myth*, 1994); RICHARD N. CURRENT, *The Lincoln Nobody Knows* (1958, reprinted 1980); DAVID HERBERT DONALD, *Lincoln's Herndon* (1948, reprinted 1988) and *Lincoln Reconsidered*, 2nd ed., enlarged (1961, reissued 1989); DON E. FEHRENBACHER, *Lincoln in Text and Context* (1987), which compiles essays on prewar politics, the Civil War, and Lincoln's changing image; GABOR S. BORITT and NORMAN O. FORNESS (eds.), *The Historian's Lincoln: Pseudohistory, Psychohistory, and History* (1988); and THOMAS J. DILORENZO, *The Real Lincoln: A New Look at Abraham Lincoln, His Agenda, and an Unnecessary War* (2002). BRUCE TAP, *Over Lincoln's Shoulder: The Committee on the Conduct of the War* (1998), offers a revealing account of the Congressional committee's interference in Lincoln's handling of the war. Lincoln's famous Second Inaugural Address is the subject of RONALD C. WHITE, JR., *Lincoln's Greatest Speech: The Second Inaugural* (2002).

Lincoln's administration is documented in PHILLIP SHAW PALU-

DAN, *The Presidency of Abraham Lincoln* (1994). Works dealing with aspects of Lincoln's statesmanship are DON E. FEHRENBACHER, *Prelude to Greatness: Lincoln in the 1850s* (1962, reissued 1970); DAVID M. POTTER, *Lincoln and His Party in the Secession Crisis* (1942, reprinted 1979), with emphasis on the period between Lincoln's election and the firing on Fort Sumter; WILLIAM B. HESSELTINE, *Lincoln and the War Governors* (1948, reissued 1972); KENNETH M. STAMPP, *And the War Came: The North and the Secession Crisis, 1860–1861* (1950, reprinted 1980); T. HARRY WILLIAMS, *Lincoln and the Radicals* (1941, reissued 1969), and *Lincoln and His Generals* (1952, reprinted 1981); HANS L. TREFOUSSE, *The Radical Republicans: Lincoln's Vanguard for Racial Justice* (1968); DAVID A. NICHOLS, *Lincoln and the Indians: Civil War Policy and Politics* (1978); GABOR S. BORITT, *Lincoln and the Economics of the American Dream* (1978); JAMES M. MCPHERSON, *Abraham Lincoln and the Second American Revolution* (1990); ROBERT W. JOHANNSEN, *Lincoln, the South, and Slavery: The Political Dimension* (1991); and MICHAEL BURLINGAME (ed.), *Lincoln Observed: Civil War Dispatches of Noah Brooks* (1998), which provides a journalist's recollections of Lincoln.

Books dealing with specific Lincoln issues are numerous. RUTH PAINTER RANDALL, *Mary Lincoln: Biography of a Marriage* (1953, reissued 1961), and *Lincoln's Sons* (1955), examine Lincoln's family life. Randall effectively refutes the views of WILLIAM H. HERNDON and JESSE W. WEIK, *Herndon's Life of Lincoln*, ed. by PAUL M. ANGLE (1930, reissued 1983). Two books about the president's wife are JEAN H. BAKER, *Mary Todd Lincoln: A Biography* (1987); and MARK E. NEELY, JR., and R. GERALD MCMURTRY, *The Insanity File: The Case of Mary Todd Lincoln* (1986, reissued 1993). An able and realistic treatment of Lincoln's legal career is JOHN J. DUFF, *A Lincoln, Prairie Lawyer* (1960). PHILIP B. KUNHARDT, JR., *A New Birth of Freedom: Lincoln at Gettysburg* (1983), focuses on aspects of his famous speech; as does GARRY WILLS, *Lincoln at Gettysburg: The Words That Remade America* (1992). HAROLD HOLZER (compiler and ed.), *Dear Mr. Lincoln: Letters to the President* (1993), assembles letters written by ordinary citizens covering all topics. HAROLD HOLZER, GABOR S. BORITT, and MARK E. NEELY, JR., *The Lincoln Image: Abraham Lincoln and the Popular Print* (1984), explores Lincoln's rise to fame through the medium of prints; and MARK S. REINHART, *Abraham Lincoln on Screen* (1999), offers a history of films and television programs that feature portrayals of Lincoln. MERRILL D. PETERSON, *Lincoln in American Memory* (1994), examines the view each succeeding generation has had toward Lincoln.

Comparisons between Lincoln and other historical figures include DAVID ZAREFSKY, *Lincoln, Douglas, and Slavery* (1990), which delves into the background of the Lincoln-Douglas debates and outlines each speaker's rhetorical methods; and WILLIAM CATTON and BRUCE CATTON, *Two Roads to Sumter* (1963, reissued 1971), which analyzes the dual roads taken by Lincoln and Jefferson Davis that led to the Civil War.

An overall view of Lincoln's assassination may be found in WILLIAM HANCHETT, *The Lincoln Murder Conspiracies* (1983). More detailed accounts of the last hours of his life include JIM BISHOP, *The Day Lincoln Was Shot* (1955, reprinted 1984); and W. EMERSON RECK, *A. Lincoln, His Last 24 Hours* (1987). A fictional account based on historical research is found in THOMAS MALLON, *Henry and Clara* (1994). (R.N.C./Ed.)

Linguistics

Linguistics is the scientific study of language. The word was first used in the middle of the 19th century to emphasize the difference between a newer approach to the study of language that was then developing and the more traditional approach of philology. The differences were and are largely matters of attitude, emphasis, and purpose. The philologist is concerned primarily with the historical development of languages as it is manifest in written texts and in the context of the associated literature and culture. The linguist, though he may be interested in written texts and in the development of languages through time, tends to give priority to spoken languages and to the problems of analyzing them as they operate at a given point in time.

The field of linguistics may be divided in terms of three dichotomies: synchronic versus diachronic, theoretical versus applied, microlinguistics versus macrolinguistics. A synchronic description of a language describes the language as it is at a given time; a diachronic description is concerned with the historical development of the language and the structural changes that have taken place in it. The goal of theoretical linguistics is the construction of a general theory of the structure of language or of a general theoretical framework for the description of languages; the

aim of applied linguistics is the application of the findings and techniques of the scientific study of language to practical tasks, especially to the elaboration of improved methods of language teaching. The terms microlinguistics and macrolinguistics are not yet well established, and they are, in fact, used here purely for convenience. The former refers to a narrower and the latter to a much broader view of the scope of linguistics. According to the microlinguistic view, languages should be analyzed for their own sake and without reference to their social function, to the manner in which they are acquired by children, to the psychological mechanisms that underlie the production and reception of speech, to the literary and the aesthetic or communicative function of language, and so on. In contrast, macrolinguistics embraces all of these aspects of language. Various areas within macrolinguistics have been given terminological recognition: psycholinguistics, sociolinguistics, anthropological linguistics, dialectology, mathematical and computational linguistics, and stylistics.

A large portion of this article is devoted to theoretical, synchronic microlinguistics, which is generally acknowledged as the central part of the subject; it will be abbreviated henceforth as theoretical linguistics. (J.Lyo.)

The article is divided into the following sections:

-
- History of linguistics 41
 - Earlier history 41
 - Non-Western traditions
 - Greek and Roman antiquity
 - The European Middle Ages
 - The Renaissance
 - The 19th century 42
 - Development of the comparative method
 - The role of analogy
 - Other 19th-century theories and development
 - The 20th century 44
 - Structuralism
 - Transformational grammar
 - Tagmemic, stratificational, and other approaches
 - Methods of synchronic linguistic analysis 45
 - Structural linguistics 45
 - Phonology
 - Morphology
 - Syntax
 - Semantics
 - Transformational-generative grammar 50
 - Harris's grammar
 - Chomsky's grammar
 - Modifications in Chomsky's grammar
 - Tagmemics 53
 - Modes of language
 - Hierarchy of levels
 - Stratificational grammar 53
 - Technical terminology
 - Interstratal relationships
 - The Prague school 54
 - Combination of structuralism and functionalism
 - Phonological contributions
 - Theory of markedness
 - Recent contributions
 - Historical (diachronic) linguistics 55
 - Linguistic change 55
 - Sound change
 - Grammatical change
 - Semantic change
 - Borrowing
 - The comparative method 56
 - Grimm's law
 - Proto-Indo-European reconstruction
 - Steps in the comparative method
 - Criticisms of the comparative method
 - Internal reconstruction
 - Language classification 58
 - Linguistics and other disciplines 59
 - Psycholinguistics 59
 - Language acquisition by children
 - Speech perception
 - Other areas of research
 - Sociolinguistics 59
 - Delineation of the field
 - Social dimensions
 - Other relationships 60
 - Anthropological linguistics
 - Computational linguistics
 - Mathematical linguistics
 - Stylistics
 - Philosophy of language
 - Applied linguistics
 - Dialectology and linguistic geography 61
 - Dialect geography 61
 - Early dialect studies
 - Dialect atlases
 - The value and applications of dialectology
 - Social dialectology 62
 - Semantics 63
 - Modern development of semantics 63
 - Positivist theory
 - Whorfian views
 - School of natural language
 - Modern grammatical influences
 - Philosophical views on meaning 64
 - Meaning and reference
 - Meaning and truth
 - Meaning and use
 - Meaning and thought
 - Meaning in linguistics 66
 - Semantics in the theory of language
 - Meaning, structure, and context
 - Lexical entries
 - Generative semantics
 - The study of writing 68
 - The scope of the study of writing 68
 - Sources of information
 - Definitions of terms
 - Grammatology 69
 - Three main approaches
 - Subdivisions of grammatology
 - Epigraphy and paleography 70
 - History of the study of writing 70
 - Studies prior to the 18th century
 - Modern Western studies
 - Bibliography 71

History of linguistics

EARLIER HISTORY

Non-Western traditions. Linguistic speculation and investigation, insofar as is known, has gone on in only a small number of societies. To the extent that Mesopotamian, Chinese, and Arabic learning dealt with grammar, their treatments were so enmeshed in the particularities of those languages and so little known to the European world until recently that they have had virtually no impact on Western linguistic tradition. Chinese linguistic and philological scholarship stretches back for more than two millennia, but the interest of those scholars was concentrated largely on phonetics, writing, and lexicography; their consideration of grammatical problems was bound up closely with the study of logic.

Certainly the most interesting non-Western grammatical tradition—and the most original and independent—is that of India, which dates back at least two and one-half millennia and which culminates with the grammar of Pāṇini, of the 5th century BC. There are three major ways in which the Sanskrit tradition has had an impact on modern linguistic scholarship. As soon as Sanskrit became known to the Western learned world the unravelling of comparative Indo-European grammar ensued and the foundations were laid for the whole 19th-century edifice of comparative philology and historical linguistics. But, for this, Sanskrit was simply a part of the data; Indian grammatical learning played almost no direct part. Nineteenth-century workers, however, recognized that the native tradition of phonetics in ancient India was vastly superior to Western knowledge; and this had important consequences for the growth of the science of phonetics in the West. Thirdly, there is in the rules or definitions (sutras) of Pāṇini a remarkably subtle and penetrating account of Sanskrit grammar. The construction of sentences, compound nouns, and the like is explained through ordered rules operating on underlying structures in a manner strikingly similar in part to modes of contemporary theory. As might be imagined, this perceptive Indian grammatical work has held great fascination for 20th-century theoretical linguists. A study of Indian logic in relation to Pāṇinian grammar alongside Aristotelian and Western logic in relation to Greek grammar and its successors could bring illuminating insights.

Whereas in ancient Chinese learning a separate field of study that might be called grammar scarcely took root, in ancient India a sophisticated version of this discipline developed early alongside the other sciences. Even though the study of Sanskrit grammar may originally have had the practical aim of keeping the sacred Vedic texts and their commentaries pure and intact, the study of grammar in India in the 1st millennium BC had already become an intellectual end in itself.

Greek and Roman antiquity. The emergence of grammatical learning in Greece is less clearly known than is sometimes implied, and the subject is more complex than is often supposed; here only the main strands can be sampled. The term *hē grammatikē technē* ("the art of letters") had two senses. It meant the study of the values of the letters and of accentuation and prosody and, in this sense, was an abstract intellectual discipline; and it also meant the skill of literacy and thus embraced applied pedagogy. This side of what was to become "grammatical" learning was distinctly applied, particular, and less exalted by comparison with other pursuits. Most of the developments associated with theoretical grammar grew out of philosophy and criticism; and in these developments a repeated duality of themes crosses and intertwines.

Much of Greek philosophy was occupied with the distinction between that which exists "by nature" and that which exists "by convention." So in language it was natural to account for words and forms as ordained by nature (by onomatopoeia—*i.e.*, by imitation of natural sounds) or as arrived at arbitrarily by a social convention. This dispute regarding the origin of language and meanings paved the way for the development of divergences between the views of the "analogists," who looked on language as possessing an essential regularity as a result of the symmetries that convention can provide, and the views of the "anomalists,"

who pointed to language's lack of regularity as one facet of the inescapable irregularities of nature. The situation was more complex, however, than this statement would suggest. For example, it seems that the anomalists among the Stoics credited the irrational quality of language precisely to the claim that language did not exactly mirror nature. In any event, the anomalist tradition in the hands of the Stoics brought grammar the benefit of their work in logic and rhetoric. This led to the distinction that, in modern theory, is made with the terms *signifiant* ("what signifies") and *signifié* ("what is signified") or, somewhat differently and more elaborately, with "expression" and "content"; and it laid the groundwork of modern theories of inflection, though by no means with the exhaustiveness and fine-grained analysis reached by the Sanskrit grammarians.

The Alexandrians, who were analogists working largely on literary criticism and text philology, completed the development of the classical Greek grammatical tradition. Dionysius Thrax, in the 2nd century BC, produced the first systematic grammar of Western tradition; it dealt only with word morphology. The study of sentence syntax was to wait for Apollonius Dyscolus, of the 2nd century AD. Dionysius called grammar "the acquaintance with [or observation of] what is uttered by poets and writers," using a word meaning a less general form of knowledge than what might be called "science." His typically Alexandrian literary goal is suggested by the headings in his work: pronunciation, poetic figurative language, difficult words, true and inner meanings of words, exposition of form-classes, literary criticism. Dionysius defined a sentence as a unit of sense or thought, but it is difficult to be sure of his precise meaning.

The Romans, who largely took over, with mild adaptations to their highly similar language, the total work of the Greeks, are important not as originators but as transmitters. Aelius Donatus, of the 4th century AD, and Priscian, an African of the 6th century, and their colleagues were slightly more systematic than their Greek models but were essentially retrospective rather than original. Up to this point a field that was at times called *ars grammatica* was a congeries of investigations, both theoretical and practical, drawn from the work and interests of literacy, scribeship, logic, epistemology, rhetoric, textual philosophy, poetics, and literary criticism. Yet modern specialists in the field still share their concerns and interests. The anomalists, who concentrated on surface irregularity and who looked then for regularities deeper down (as the Stoics sought them in logic) bear a resemblance to contemporary scholars of the transformationalist school. And the philological analogists with their regularizing surface segmentation show striking kinship of spirit with the modern school of structural (or taxonomic or glossematic) grammatical theorists.

The European Middle Ages. It is possible that developments in grammar during the Middle Ages constitute one of the most misunderstood areas of the field of linguistics. It is difficult to relate this period coherently to other periods and to modern concerns because surprisingly little is accessible and certain, let alone analyzed with sophistication. In the early 1970s the majority of the known grammatical treatises had not yet been made available in full to modern scholarship, so that not even their true extent could be classified with confidence. These works must be analyzed and studied in the light of medieval learning, especially the learning of the schools of philosophy then current, in order to understand their true value and place.

The field of linguistics has almost completely neglected the achievements of this period. Students of grammar have tended to see as high points in their field the achievements of the Greeks, the Renaissance growth and "rediscovery" of learning (which led directly to modern school traditions), the contemporary flowering of theoretical study (men usually find their own age important and fascinating), and, in recent decades, the astonishing monument of Pāṇini. Many linguists have found uncongenial the combination of medieval Latin learning and premodern philosophy. Yet medieval scholars might reasonably be expected to have bequeathed to modern scholarship the fruits of more than ordinarily refined perceptions of a certain order. These scholars used, wrote in, and studied

Sanskrit
grammar

"Analogists" and
"anomalists"

Use of
Latin by
medieval
scholars

Latin, a language that, though not their native tongue, was one in which they were very much at home; such scholars in groups must often have represented a highly varied linguistic background.

Some of the medieval treatises continue the tradition of grammars of late antiquity; so there are versions based on Donatus and Priscian, often with less incorporation of the classical poets and writers. Another genre of writing involves simultaneous consideration of grammatical distinctions and scholastic logic; modern linguists are probably inadequately trained to deal with these writings.

Certainly the most obviously interesting theorizing to be found in this period is contained in the "speculative grammar" of the *modistae*, who were so called because the titles of their works were often phrased *De modis significandi tractatus* ("Treatise Concerning the Modes of Signifying"). For the development of the Western grammatical tradition, work of this genre was the second great milestone after the crystallization of Greek thought with the Stoics and Alexandrians. The scholastic philosophers were occupied with relating words and things—*i.e.*, the structure of sentences with the nature of the real world—hence their preoccupation with signification. The aim of the grammarians was to explore how a word (an element of language) matched things apprehended by the mind and how it signified reality. Since a word cannot signify the nature of reality directly, it must stand for the thing signified in one of its modes or properties; it is this discrimination of modes that the study of categories and parts of speech is all about. Thus the study of sentences should lead one to the nature of reality by way of the modes of signifying.

The *modistae* did not innovate in discriminating categories and parts of speech; they accepted those that had come down from the Greeks through Donatus and Priscian. The great contribution of these grammarians, who flourished between the mid-13th and mid-14th century, was their insistence on a grammar to explicate the distinctions found by their forerunners in the languages known to them. Whether they made the best choice in selecting logic, metaphysics, and epistemology (as they knew them) as the fields to be included with grammar as a basis for the grand account of universal knowledge is less important than the breadth of their conception of the place of grammar. Before the *modistae*, grammar had not been viewed as a separate discipline but had been considered in conjunction with other studies or skills (such as criticism, preservation of valued texts, foreign-language learning). The Greek view of grammar was rather narrow and fragmented; the Roman view was largely technical. The speculative medieval grammarians (who dealt with language as a *speculum*, "mirror" of reality) inquired into the fundamentals underlying language and grammar. They wondered whether grammarians or philosophers discovered grammar, whether grammar was the same for all languages, what the fundamental topic of grammar was, and what the basic and irreducible grammatical primes are. Signification was reached by imposition of words on things; *i.e.*, the sign was arbitrary. Those questions sound remarkably like current issues of linguistics, which serves to illustrate how slow and repetitious progress in the field is. While the *modistae* accepted, by modern standards, a restrictive set of categories, the acumen and sweep they brought to their task resulted in numerous subtle and fresh syntactic observations. A thorough study of the medieval period would greatly enrich the discussion of current questions.

The Renaissance. It is customary to think of the Renaissance as a time of great flowering. There is no doubt that linguistic and philological developments of this period are interesting and significant. Two new sets of data that modern linguists tend to take for granted became available to grammarians during this period: (1) the newly recognized vernacular languages of Europe, for the protection and cultivation of which there subsequently arose national academies and learned institutions that live down to the present day; and (2) the exotic languages of Africa, the Orient, the New World, and, later, of Siberia, Inner Asia, Papua, Oceania, the Arctic, and Australia, which the

voyages of discovery opened up. Earlier, the only non-Indo-European grammar at all widely accessible was that of the Hebrews (and to some extent Arabic); and Semitic in fact shares many categories with Indo-European in its grammar. Indeed, for many of the exotic languages scholarship barely passed beyond the most rudimentary initial collection of word lists; grammatical analysis was scarcely approached.

In the field of grammar, the Renaissance did not produce notable innovation or advance. Generally speaking, there was a strong rejection of speculative grammar and a relatively uncritical resumption of late Roman views (as stated by Priscian). This was somewhat understandable in the case of Latin or Greek grammars, since here the task was less evidently that of intellectual inquiry and more that of the schools, with the practical aim of gaining access to the newly discovered ancients. But, aside from the fact that, beginning in the 15th century, serious grammars of European vernaculars were actually written, it is only in particular cases and for specific details (*e.g.*, a mild alteration in the number of parts of speech or cases of nouns) that real departures from Roman grammar can be noted. Likewise, until the end of the 19th century, grammars of the exotic languages, written largely by missionaries and traders, were cast almost entirely in the Roman model, to which the Renaissance had added a limited medieval syntactic ingredient.

From time to time a degree of boldness may be seen in France: Petrus Ramus, a 16th-century logician, worked within a taxonomic framework of the surface shapes of words and inflections, such work entailing some of the attendant trivialities that modern linguistics has experienced (*e.g.*, by dividing up Latin nouns on the basis of equivalence of syllable count among their case forms). In the 17th century, members of Solitaires (a group of hermits who lived in the deserted abbey of Port-Royal in France) produced a grammar that has exerted noteworthy continuing influence, even in contemporary theoretical discussion. Drawing their basic view from scholastic logic as modified by rationalism, these people aimed to produce a philosophical grammar that would capture what was common to the grammars of languages—a general grammar, but not aprioristically universalist. This grammar has attracted recent attention because it employs certain syntactic formulations that resemble in detail contemporary transformational rules, which formulate the relationship between the various elements of a sentence.

Roughly from the 15th century to World War II, however, the version of grammar available to the Western public (together with its colonial expansion) remained basically that of Priscian with only occasional and subsidiary modifications, and the knowledge of new languages brought only minor adjustments to the serious study of grammar. As education has become more broadly disseminated throughout society by the schools, attention has shifted from theoretical or technical grammar as an intellectual preoccupation to prescriptive grammar suited to pedagogical purposes, which started with Renaissance vernacular nationalism. Grammar increasingly parted company with its older fellow disciplines within philosophy as they moved over to the domain known as natural science, and technical academic grammatical study has increasingly become involved with issues represented by empiricism versus rationalism and their successor manifestations on the academic scene.

Nearly down to the present day, the grammar of the schools has had only tangential connections with the studies pursued by professional linguists; for most people prescriptive grammar has become synonymous with "grammar," and the prevailing view held by educated people regards grammar as an item of folk knowledge open to speculation by all, and in nowise a formal science requiring adequate preparation such as is assumed for chemistry. (E.P.H./Ed.)

THE 19TH CENTURY

Development of the comparative method. It is generally agreed that the most outstanding achievement of linguistic scholarship in the 19th century was the development of

The
modistae

The Port-
Royal
school

New
sources of
data

Discovery of the Indo-European language family

the comparative method, which comprised a set of principles whereby languages could be systematically compared with respect to their sound systems, grammatical structure, and vocabulary and shown to be "genealogically" related. As French, Italian, Portuguese, Romanian, Spanish, and the other Romance languages had evolved from Latin, so Latin, Greek, and Sanskrit as well as the Celtic, Germanic, and Slavic languages and many other languages of Europe and Asia had evolved from some earlier language, to which the name Indo-European or Proto-Indo-European is now customarily applied. That all the Romance languages were descended from Latin and thus constituted one "family" had been known for centuries; but the existence of the Indo-European family of languages and the nature of their genealogical relationship was first demonstrated by the 19th-century comparative philologists. (The term philology in this context is not restricted to the study of literary languages.)

The main impetus for the development of comparative philology came toward the end of the 18th century, when it was discovered that Sanskrit bore a number of striking resemblances to Greek and Latin. An English orientalist, Sir William Jones, though he was not the first to observe these resemblances, is generally given the credit for bringing them to the attention of the scholarly world and putting forward the hypothesis, in 1786, that all three languages must have "sprung from some common source, which perhaps no longer exists." By this time, a number of texts and glossaries of the older Germanic languages (Gothic, Old High German, and Old Norse) had been published, and Jones realized that Germanic as well as Old Persian and perhaps Celtic had evolved from the same "common source." The next important step came in 1822, when the German scholar Jacob Grimm, following the Danish linguist Rasmus Rask (whose work, being written in Danish, was less accessible to most European scholars), pointed out in the second edition of his comparative grammar of Germanic that there were a number of systematic correspondences between the sounds of Germanic and the sounds of Greek, Latin, and Sanskrit in related words. Grimm noted, for example, that where Gothic (the oldest surviving Germanic language) had an *f*, Latin, Greek, and Sanskrit frequently had a *p* (e.g., Gothic *fotus*, Latin *pedis*, Greek *podós*, Sanskrit *padás*, all meaning "foot"); when Gothic had a *p*, the non-Germanic languages had a *b*; when Gothic had a *b*, the non-Germanic languages had what Grimm called an "aspirate" (Latin *f*, Greek *ph*, Sanskrit *bh*). In order to account for these correspondences he postulated a cyclical "soundshift" (*Lautverschiebung*) in the prehistory of Germanic, in which the original "aspirates" became voiced unaspirated stops (*bh* became *b*, etc.), the original voiced unaspirated stops became voiceless (*b* became *p*, etc.), and the original voiceless (unaspirated) stops became "aspirates" (*p* became *f*). Grimm's term, "aspirate," it will be noted, covered such phonetically distinct categories as aspirated stops (*bh*, *ph*), produced with an accompanying audible puff of breath, and fricatives (*f*), produced with audible friction as a result of incomplete closure in the vocal tract.

The Neogrammarians and sound laws

In the work of the next 50 years the idea of sound change was made more precise, and, in the 1870s, a group of scholars known collectively as the *Junggrammatiker* ("young grammarians," or Neogrammarians) put forward the thesis that all changes in the sound system of a language as it developed through time were subject to the operation of regular sound laws. Though the thesis that sound laws were absolutely regular in their operation (unless they were inhibited in particular instances by the influence of analogy) was at first regarded as most controversial, by the end of the 19th century it was quite generally accepted and had become the cornerstone of the comparative method. Using the principle of regular sound change, scholars were able to reconstruct "ancestral" common forms from which the later forms found in particular languages could be derived. By convention, such reconstructed forms are marked in the literature with an asterisk. Thus, from the reconstructed Proto-Indo-European word for "ten," **dek̑m*, it was possible to derive Sanskrit *daśa*, Greek *déka*, Latin *decem*, and Gothic

taihun by postulating a number of different sound laws that operated independently in the different branches of the Indo-European family. The question of sound change is dealt with in greater detail in the section entitled *Historical linguistics*.

The role of analogy. Analogy has been mentioned in connection with its inhibition of the regular operation of sound laws in particular word forms. This was how the Neogrammarians thought of it. In the course of the 20th century, however, it has come to be recognized that analogy, taken in its most general sense, plays a far more important role in the development of languages than simply that of sporadically preventing what would otherwise be a completely regular transformation of the sound system of a language. When a child learns to speak he tends to regularize the anomalous, or irregular, forms by analogy with the more regular and productive patterns of formation in the language; e.g., he will tend to say "comed" rather than "came," "dived" rather than "dove," and so on, just as he will say "talked," "loved," and so forth. The fact that the child does this is evidence that he has learned or is learning the regularities or rules of his language. He will go on to "unlearn" some of the analogical forms and substitute for them the anomalous forms current in the speech of the previous generation. But in some cases, he will keep a "new" analogical form (e.g., "dived" rather than "dove"), and this may then become the recognized and accepted form.

Other 19th-century theories and development. *Inner and outer form.* One of the most original, if not one of the most immediately influential, linguists of the 19th century was the learned Prussian statesman, Wilhelm von Humboldt (died 1835). His interests, unlike those of most of his contemporaries, were not exclusively historical. Following the German philosopher Johann Gottfried von Herder (1744–1803), he stressed the connection between national languages and national character: this was but a commonplace of romanticism. More original was Humboldt's theory of "inner" and "outer" form in language. The outer form of language was the raw material (the sounds) from which different languages were fashioned; the inner form was the pattern, or structure, of grammar and meaning that was imposed upon this raw material and differentiated one language from another. This "structural" conception of language was to become dominant, for a time at least, in many of the major centres of linguistics by the middle of the 20th century. Another of Humboldt's ideas was that language was something dynamic, rather than static, and was an activity itself rather than the product of activity. A language was not a set of actual utterances produced by speakers but the underlying principles or rules that made it possible for speakers to produce such utterances and, moreover, an unlimited number of them. This idea was taken up by a German philologist, Heymann Steintal, and, what is more important, by the physiologist and psychologist Wilhelm Wundt, and thus influenced late 19th- and early 20th-century theories of the psychology of language. Its influence, like that of the distinction of inner and outer form, can also be seen in the thought of Ferdinand de Saussure, a Swiss linguist. But its full implications were probably not perceived and made precise until the middle of the 20th century, when the U.S. linguist Noam Chomsky re-emphasized it and made it one of the basic notions of generative grammar (see below).

Phonetics and dialectology. Many other interesting and important developments occurred in 19th-century linguistic research, among them work in the areas of phonetics and dialectology. Research in both these fields was promoted by the Neogrammarians' concern with sound change and by their insistence that prehistoric developments in languages were of the same kind as developments taking place in the languages and dialects currently spoken. The development of phonetics in the West was also strongly influenced at this period, as were many of the details of the more philological analysis of the Indo-European languages, by the discovery of the works of the Indian grammarians who, from the time of the Sanskrit grammarian Pāṇini (5th or 6th century BC), if not before, had arrived at a much more comprehensive and scientific

Ideas of Wilhelm von Humboldt

Development of phonetics

theory of phonetics, phonology, and morphology than anything achieved in the West until the modern period.

THE 20TH CENTURY

Structuralism. The term structuralism has been used as a slogan and rallying cry by a number of different schools of linguistics, and it is necessary to realize that it has somewhat different implications according to the context in which it is employed. It is convenient to draw first a broad distinction between European and American structuralism and, then, to treat them separately.

Structural linguistics in Europe. Structural linguistics in Europe is generally said to have begun in 1916 with the posthumous publication of the *Cours de Linguistique Générale* (*Course in General Linguistics*) of Ferdinand de Saussure. Much of what is now considered as Saussurean can be seen, though less clearly, in the earlier work of Humboldt, and the general structural principles that Saussure was to develop with respect to synchronic linguistics in the *Cours* had been applied almost 40 years before (1879) by Saussure himself in a reconstruction of the Indo-European vowel system. The full significance of the work was not appreciated at the time. Saussure's structuralism can be summed up in two dichotomies (which jointly cover what Humboldt referred to in terms of his own distinction of inner and outer form): (1) *langue* versus *parole* and (2) form versus substance. By *langue*, best translated in its technical Saussurean sense as language system, is meant the totality of regularities and patterns of formation that underlie the utterances of a language; by *parole*, which can be translated as language behaviour, is meant the actual utterances themselves. Just as two performances of a piece of music given by different orchestras on different occasions will differ in a variety of details and yet be identifiable as performances of the same piece, so two utterances may differ in various ways and yet be recognized as instances, in some sense, of the same utterance. What the two musical performances and the two utterances have in common is an identity of form, and this form, or structure, or pattern, is in principle independent of the substance, or "raw material," upon which it is imposed. "Structuralism," in the European sense then, refers to the view that there is an abstract relational structure that underlies and is to be distinguished from actual utterances—a system underlying actual behaviour—and that this is the primary object of study for the linguist.

Two important points arise here: first, that the structural approach is not in principle restricted to synchronic linguistics; second, that the study of meaning, as well as the study of phonology and grammar, can be structural in orientation. In both cases "structuralism" is opposed to "atomism" in the European literature. It was Saussure who drew the terminological distinction between synchronic and diachronic linguistics in the *Cours*: despite the undoubtedly structural orientation of his own early work in the historical and comparative field, he maintained that, whereas synchronic linguistics should deal with the structure of a language system at a given point in time, diachronic linguistics should be concerned with the historical development of isolated elements—it should be atomistic. Whatever the reasons that led Saussure to take this rather paradoxical view, his teaching on this point was not generally accepted, and scholars soon began to apply structural concepts to the diachronic study of languages. The most important of the various schools of structural linguistics to be found in Europe in the first half of the 20th century have included the Prague school, most notably represented by Nikolay Sergeevich Trubetskoy (died 1938) and Roman Jakobson (born 1896), both Russian émigrés, and the Copenhagen (or glossematic) school, centred around Louis Hjelmslev (died 1965). John Rupert Firth (died 1960) and his followers, sometimes referred to as the London school, were less Saussurean in their approach, but, in a general sense of the term, their approach may also be described appropriately as structural linguistics.

Structural linguistics in America. American and European structuralism shared a number of features. In insisting upon the necessity of treating each language as a more or less coherent and integrated system, both European and

American linguists of this period tended to emphasize, if not to exaggerate, the structural uniqueness of individual languages. There was especially good reason to take this point of view given the conditions in which American linguistics developed from the end of the 19th century. There were hundreds of indigenous American Indian languages that had never been previously described. Many of these were spoken by only a handful of speakers and, if they were not recorded before they became extinct, would be permanently inaccessible. Under these circumstances, such linguists as Franz Boas (died 1942) were less concerned with the construction of a general theory of the structure of human language than they were with prescribing sound methodological principles for the analysis of unfamiliar languages. They were also fearful that the description of these languages would be distorted by analyzing them in terms of categories derived from the analysis of the more familiar Indo-European languages.

After Boas, the two most influential American linguists were Edward Sapir (died 1939) and Leonard Bloomfield (died 1949). Like his teacher Boas, Sapir was equally at home in anthropology and linguistics, the alliance of which disciplines has endured to the present day in many American universities. Boas and Sapir were both attracted by the Humboldtian view of the relationship between language and thought, but it was left to one of Sapir's pupils, Benjamin Lee Whorf, to present it in a sufficiently challenging form to attract widespread scholarly attention. Since the republication of Whorf's more important papers in 1956, the thesis that language determines perception and thought has come to be known as the Whorfian hypothesis.

Sapir's work has always held an attraction for the more anthropologically inclined American linguists. But it was Bloomfield who prepared the way for the later phase of what is now thought of as the most distinctive manifestation of American "structuralism." When he published his first book in 1914, Bloomfield was strongly influenced by Wundt's psychology of language. In 1933, however, he published a drastically revised and expanded version with the new title *Language*; this book dominated the field for the next 30 years. In it Bloomfield explicitly adopted a behaviouristic approach to the study of language, eschewing in the name of scientific objectivity all reference to mental or conceptual categories. Of particular consequence was his adoption of the behaviouristic theory of semantics according to which meaning is simply the relationship between a stimulus and a verbal response. Because science was still a long way from being able to give a comprehensive account of most stimuli, no significant or interesting results could be expected from the study of meaning for some considerable time, and it was preferable, as far as possible, to avoid basing the grammatical analysis of a language on semantic considerations. Bloomfield's followers pushed even further the attempt to develop methods of linguistic analysis that were not based on meaning. One of the most characteristic features of "post-Bloomfieldian" American structuralism, then, was its almost complete neglect of semantics.

Another characteristic feature, one that was to be much criticized by Chomsky, was its attempt to formulate a set of "discovery procedures"—procedures that could be applied more or less mechanically to texts and could be guaranteed to yield an appropriate phonological and grammatical description of the language of the texts. Structuralism, in this narrower sense of the term, is represented, with differences of emphasis or detail, in the major American textbooks published during the 1950s.

Transformational grammar. The most significant development in linguistic theory and research in recent years was the rise of generative grammar, and, more especially, of transformational-generative grammar, or transformational grammar, as it came to be known. Two versions of transformational grammar were put forward in the mid-1950s, the first by Zellig S. Harris and the second by Noam Chomsky, his pupil. It is Chomsky's system that has attracted the most attention so far. As first presented by Chomsky in *Syntactic Structures* (1957), transformational grammar can be seen partly as a reaction against

The work of Sapir, Bloomfield, and Whorf

European sense of structuralism

The work of Harris and Chomsky

post-Bloomfieldian structuralism and partly as a continuation of it. What Chomsky reacted against most strongly was the post-Bloomfieldian concern with discovery procedures. In his opinion, linguistics should set itself the more modest and more realistic goal of formulating criteria for evaluating alternative descriptions of a language without regard to the question of how these descriptions had been arrived at. The statements made by linguists in describing a language should, however, be cast within the framework of a far more precise theory of grammar than had hitherto been the case, and this theory should be formalized in terms of modern mathematical notions. Within a few years, Chomsky had broken with the post-Bloomfieldians on a number of other points also. He had adopted what he called a "mentalist" theory of language, by which term he implied that the linguist should be concerned with the speaker's creative linguistic competence and not his performance, the actual utterances produced. He had challenged the post-Bloomfieldian concept of the phoneme (see below), which many scholars regarded as the most solid and enduring result of the previous generation's work. And he had challenged the structuralists' insistence upon the uniqueness of every language, claiming instead that all languages were, to a considerable degree, cut to the same pattern—they shared a certain number of formal and substantive universals.

Tagmemic, stratificational, and other approaches. The effect of Chomsky's ideas has been phenomenal. It is hardly an exaggeration to say that there is no major theoretical issue in linguistics today that is debated in terms other than those in which he has chosen to define it, and every school of linguistics tends to define its position in relation to his. Among the rival schools are tagmemics, stratificational grammar, and the Prague school. Tagmemic is the system of linguistic analysis developed by the U.S. linguist Kenneth L. Pike and his associates in connection with their work as Bible translators. Its foundations were laid during the 1950s, when Pike differed from the post-Bloomfieldian structuralists on a number of principles, and it has been further elaborated since then. Tagmemic analysis has been used for analyzing a great many previously unrecorded languages, especially in Central and South America and in West Africa. Stratificational grammar, developed by a U.S. linguist, Sydney M. Lamb, has been seen by some linguists as an alternative to transformational grammar. Not yet fully expounded or widely exemplified in the analysis of different languages, stratificational grammar is perhaps best characterized as a radical modification of post-Bloomfieldian linguistics, but it has many features that link it with European structuralism. The Prague school has been mentioned above for its importance in the period immediately following the publication of Saussure's *Cours*. Many of its characteristic ideas (in particular, the notion of distinctive features in phonology) have been taken up by other schools. But there has been further development in Prague of the functional approach to syntax (see below). The work of M.A.K. Halliday in England derived much of its original inspiration from Firth (above), but Halliday provided a more systematic and comprehensive theory of the structure of language than Firth had, and it has been quite extensively illustrated.

Methods of synchronic linguistic analysis

STRUCTURAL LINGUISTICS

This section is concerned mainly with a version of structuralism (which may also be called descriptive linguistics) developed by scholars working in a post-Bloomfieldian tradition.

Phonology. With the great progress made in phonetics in the late 19th century, it had become clear that the question whether two speech sounds were the same or not was more complex than might appear at first sight. Two utterances of what was taken to be the same word might differ quite perceptibly from one occasion of utterance to the next. Some of this variation could be attributed to a difference of dialect or accent and is of no concern here. But even two utterances of the same word by the

same speaker might vary from one occasion to the next. Variation of this kind, though it is generally less obvious and would normally pass unnoticed, is often clear enough to the trained phonetician and is measurable instrumentally. It is known that the "same" word is being uttered, even if the physical signal produced is variable, in part, because the different pronunciations of the same word will cluster around some acoustically identifiable norm. But this is not the whole answer, because it is actually impossible to determine norms of pronunciation in purely acoustic terms. Once it has been decided what counts as "sameness" of sound from the linguistic point of view, the permissible range of variation for particular sounds in particular contexts can be measured, and, within certain limits, the acoustic cues for the identification of utterances as "the same" can be determined.

What is at issue is the difference between phonetic and phonological (or phonemic) identity, and for these purposes it will be sufficient to define phonetic identity in terms solely of acoustic "sameness." Absolute phonetic identity is a theoretical ideal never fully realized. From a purely phonetic point of view, sounds are more or less similar, rather than absolutely the same or absolutely different. Speech sounds considered as units of phonetic analysis in this article are called phones, and, following the normal convention, are represented by enclosing the appropriate alphabetic symbol in square brackets. Thus [p] will refer to a *p* sound (*i.e.*, what is described more technically as a voiceless, bilabial stop); and [pit] will refer to a complex of three phones—a *p* sound, followed by an *i* sound, followed by a *t* sound. A phonetic transcription may be relatively broad (omitting much of the acoustic detail) or relatively narrow (putting in rather more of the detail), according to the purpose for which it is intended. A very broad transcription will be used in this article except when finer phonetic differences must be shown.

Phonological, or phonemic, identity was referred to above as "sameness of sound from the linguistic point of view." Considered as phonological units—*i.e.*, from the point of view of their function in the language—sounds are described as phonemes and are distinguished from phones by enclosing their appropriate symbol (normally, but not necessarily, an alphabetic one) between two slash marks. Thus /p/ refers to a phoneme that may be realized on different occasions of utterance or in different contexts by a variety of more or less different phones. Phonological identity, unlike phonetic similarity, is absolute: two phonemes are either the same or different, they cannot be more or less similar. For example, the English words "bit" and "pit" differ phonemically in that the first has the phoneme /b/ and the second has the phoneme /p/ in initial position. As the words are normally pronounced, the phonetic realization of /b/ will differ from the phonetic realization of /p/ in a number of different ways: it will be at least partially voiced (*i.e.*, there will be some vibration of the vocal cords), it will be without aspiration (*i.e.*, there will be no accompanying slight puff of air, as there will be in the case of the phone realizing /p/), and it will be pronounced with less muscular tension. It is possible to vary any one or all of these contributory differences, making the phones in question more or less similar, and it is possible to reduce the phonetic differences to the point that the hearer cannot be certain which word, "bit" or "pit," has been uttered. But it must be either one or the other; there is no word with an initial sound formed in the same manner as /p/ or /b/ that is halfway between the two. This is what is meant by saying that phonemes are absolutely distinct from one another—they are discrete rather than continuously variable.

How it is known whether two phones realize the same phoneme or not is dealt with differently by different schools of linguists. The "orthodox" post-Bloomfieldian school regards the first criterion to be phonetic similarity. Two phones are not said to realize the same phoneme unless they are sufficiently similar. What is meant by "sufficiently similar" is rather vague, but it must be granted that for every phoneme there is a permissible range of variation in the phones that realize it. As far as occurrence in the same context goes, there are no serious problems.

Phones:
speech
sounds

Prague
school

Classifying
phones or
phonemes

More critical is the question of whether two phones occurring in different contexts can be said to realize the same phoneme or not. To take a standard example from English: the phone that occurs at the beginning of the word "pit" differs from the phone that occurs after the initial /s/ of "spit." The "p sound" occurring after the /s/ is unaspirated (*i.e.*, it is pronounced without any accompanying slight puff of air). The aspirated and unaspirated "p sounds" may be symbolized rather more narrowly as [p^h] and [p] respectively. The question then is whether [p^h] and [p] realize the same phoneme /p/ or whether each realizes a different phoneme. They satisfy the criterion of phonetic similarity, but this, though a necessary condition of phonemic identity, is not a sufficient one.

The next question is whether there is any pair of words in which the two phones are in minimal contrast (or opposition); that is, whether there is any context in English in which the occurrence of the one rather than the other has the effect of distinguishing two or more words (in the way that [p^h] versus [b] distinguishes the so-called minimal pairs "pit" and "bit," "pan" and "ban," and so on). If there is, it can be said that, despite their phonetic similarity, the two phones realize (or "belong to") different phonemes—that the difference between them is phonemic. If there is no context in which the two phones are in contrast (or opposition) in this sense, it can be said that they are variants of the same phoneme—that the difference between them is nonphonemic. Thus, the difference between [p^h] and [p] in English is nonphonemic; the two sounds realize, or belong to, the same phoneme, namely /p/. In several other languages—*e.g.*, Hindi—the contrast between such sounds as [p^h] and [p] is phonemic, however. The question is rather more complicated than it has been represented here. In particular, it should be noted that [p] is phonetically similar to [b] as well as to [p^h] and that, although [p^h] and [b] are in contrast, [p] and [b] are not. It would thus be possible to regard [p] and [b] as variants of the same phoneme. Most linguists, however, have taken the alternative view, assigning [p] to the same phoneme as [p^h]. Here it will suffice to note that the criteria of phonetic similarity and lack of contrast do not always uniquely determine the assignment of phones to phonemes. Various supplementary criteria may then be invoked.

Phones that can occur and do not contrast in the same context are said to be in free variation in that context, and, as has been shown, there is a permissible range of variation for the phonetic realization of all phonemes. More important than free variation in the same context, however, is systematically determined variation according to the context in which a given phoneme occurs. To return to the example used above: [p] and [p^h], though they do not contrast, are not in free variation either. Each of them has its own characteristic positions of occurrence, and neither occurs, in normal English pronunciation, in any context characteristic for the other (*e.g.*, only [p^h] occurs at the beginning of a word, and only [p] occurs after s). This is expressed by saying that they are in complementary distribution. (The distribution of an element is the whole range of contexts in which it can occur.) Granted that [p] and [p^h] are variants of the same phoneme /p/, it can be said that they are contextually, or positionally, determined variants of it. To use the technical term, they are allophones of /p/. The allophones of a phoneme, then, are its contextually determined variants and they are in complementary distribution.

The post-Bloomfieldians made the assignment of phones to phonemes subject to what is now generally referred to as the principle of bi-uniqueness. The phonemic specification of a word or utterance was held to determine uniquely its phonetic realization (except for free variation), and, conversely, the phonetic description of a word or utterance was held to determine uniquely its phonemic analysis. Thus, if two words or utterances are pronounced alike, then they must receive the same phonemic description; conversely, two words or utterances that have been given the same phonemic analysis must be pronounced alike. The principle of bi-uniqueness was also held to imply that, if a given phone was assigned to a particular

phoneme in one position of occurrence, then it must be assigned to the same phoneme in all its other positions of occurrence; it could not be the allophone of one phoneme in one context and of another phoneme in other contexts.

A second important principle of the post-Bloomfieldian approach was its insistence that phonemic analysis should be carried out prior to and independently of grammatical analysis. Neither this principle nor that of bi-uniqueness was at all widely accepted outside the post-Bloomfieldian school, and they have been abandoned by the generative phonologists (see below).

Phonemes of the kind referred to so far are segmental; they are realized by consonantal or vocalic (vowel) segments of words, and they can be said to occur in a certain order relative to one another. For example, in the phonemic representation of the word "bit," the phoneme /b/ precedes /i/, which precedes /t/. But nonsegmental, or suprasegmental, aspects of the phonemic realization of words and utterances may also be functional in a language. In English, for example, the noun "import" differs from the verb "import" in that the former is accented on the first and the latter on the second syllable. This is called a stress accent: the accented syllable is pronounced with greater force or intensity. Many other languages distinguish words suprasegmentally by tone. For example, in Mandarin Chinese the words *hào* "day" and *hǎo* "good" are distinguished from one another in that the first has a falling tone and the second a falling-rising tone; these are realized, respectively, as (1) a fall in the pitch of the syllable from high to low and (2) a change in the pitch of the syllable from medium to low and back to medium. Stress and tone are suprasegmental in the sense that they are "superimposed" upon the sequence of segmental phonemes. The term tone is conventionally restricted by linguists to phonologically relevant variations of pitch at the level of words. Intonation, which is found in all languages, is the variation in the pitch contour or pitch pattern of whole utterances, of the kind that distinguishes (either of itself or in combination with some other difference) statements from questions or indicates the mood or attitude of the speaker (as hesitant, surprised, angry, and so forth). Stress, tone, and intonation do not exhaust the phonologically relevant suprasegmental features found in various languages, but they are among the most important.

A complete phonological description of a language includes all the segmental phonemes and specifies which allophones occur in which contexts. It also indicates which sequences of phonemes are possible in the language and which are not: it will indicate, for example, that the sequences /bl/ and /br/ are possible at the beginning of English words but not /bn/ or /bm/. A phonological description also identifies and states the distribution of the suprasegmental features. Just how this is to be done, however, has been rather more controversial in the post-Bloomfieldian tradition. Differences between the post-Bloomfieldian approach to phonology and approaches characteristic of other schools of structural linguistics will be treated below.

Morphology. The grammatical description of many, if not all, languages is conveniently divided into two complementary sections: morphology and syntax. The relationship between them, as generally stated, is as follows: morphology accounts for the internal structure of words, and syntax describes how words are combined to form phrases, clauses, and sentences.

There are many words in English that are fairly obviously analyzable into smaller grammatical units. For example, the word "unacceptability" can be divided into *un-*, *accept-*, *abil-*, and *-ity* (*abil-* being a variant of *-able*). Of these, at least three are minimal grammatical units, in the sense that they cannot be analyzed into yet smaller grammatical units—*un-*, *abil-*, and *ity*. The status of *accept-* from this point of view, is somewhat uncertain. Given the existence of such forms as *accede* and *accuse*, on the one hand, and of *except*, *excede*, and *excuse*, on the other, one might be inclined to analyze *accept* into *ac-* (which might subsequently be recognized as a variant of *ad-*) and *-cept*. The question is left open. Minimal grammatical units like *un-*, *abil-*, and *-ity* are what Bloomfield called mor-

Supra-
segmental
phonemes

Free variation and complementary distribution

Distinction between morphology and syntax

phemes; he defined them in terms of the "partial phonetic-semantic resemblance" holding within sets of words. For example, "unacceptable," "untrue," and "ungracious" are phonetically (or, phonologically) similar as far as the first syllable is concerned and are similar in meaning in that each of them is negative by contrast with a corresponding positive adjective ("acceptable," "true," "gracious"). This "partial phonetic-semantic resemblance" is accounted for by noting that the words in question contain the same morpheme (namely, *un-*) and that this morpheme has a certain phonological form and a certain meaning.

Bloomfield's definition of the morpheme in terms of "partial phonetic-semantic resemblance" was considerably modified and, eventually, abandoned entirely by some of his followers. Whereas Bloomfield took the morpheme to be an actual segment of a word, others defined it as being a purely abstract unit, and the term morph was introduced to refer to the actual word segments. The distinction between morpheme and morph (which is, in certain respects, parallel to the distinction between phoneme and phone) may be explained by means of an example. If a morpheme in English is posited with the function of accounting for the grammatical difference between singular and plural nouns, it may be symbolized by enclosing the term plural within brace brackets. Now the morpheme [plural] is represented in a number of different ways. Most plural nouns in English differ from the corresponding singular forms in that they have an additional final segment. In the written forms of these words, it is either *-s* or *-es* (e.g., "cat" : "cats"; "dog" : "dogs"; "fish" : "fishes"). The word segments written *-s* or *-es* are morphs. So also is the word segment written *-en* in "oxen." All these morphs represent the same morpheme. But there are other plural nouns in English that differ from the corresponding singular forms in other ways (e.g., "mouse" : "mice"; "criterion" : "criteria"; and so on) or not at all (e.g., "this sheep" : "these sheep"). Within the post-Bloomfieldian framework no very satisfactory account of the formation of these nouns could be given. But it was clear that they contained (in some sense) the same morpheme as the more regular plurals.

Morphs that are in complementary distribution and represent the same morpheme are said to be allomorphs of that morpheme. For example, the regular plurals of English nouns are formed by adding one of three morphs on to the form of the singular: /s/, /z/, or /ɪz/ (in the corresponding written forms both /s/ and /z/ are written *-s* and /ɪz/ is written *-es*). Their distribution is determined by the following principle: if the morph to which they are to be added ends in a "sibilant" sound (e.g., *s*, *z*, *sh*, *ch*), then the syllabic allomorph /ɪz/ is selected (e.g., *fishes* /fɪʃ-ɪz/, *matches* /mætʃ-ɪz/); otherwise the nonsyllabic allomorphs are selected, the voiceless allomorph /s/ with morphs ending in a voiceless consonant (e.g., *cat-s* /kæt-s/) and the voiced allomorph /z/ with morphs ending in a vowel or voiced consonant (e.g., *flea-s* /fli-z/, *dog-s* /dɒg-z/). These three allomorphs, it will be evident, are in complementary distribution, and the alternation between them is determined by the phonological structure of the preceding morph. Thus the choice is phonologically conditioned.

Very similar is the alternation between the three principal allomorphs of the past participle ending, /ɪd/, /t/, and /d/, all of which correspond to the *-ed* of the written forms. If the preceding morph ends with /t/ or /d/, then the syllabic allomorph /ɪd/ is selected (e.g., *wait-ed* /weɪt-ɪd/). Otherwise, if the preceding morph ends with a voiceless consonant, one of the nonsyllabic allomorphs is selected—the voiceless allomorph /t/ when the preceding morph ends with a voiceless consonant (e.g., *pack-ed* /pæk-t/) and the voiced allomorph /d/ when the preceding morph ends with a vowel or voiced consonant (e.g., *row-ed* /rou-d/; *tame-d* /teɪm-d/). This is another instance of phonological conditioning. Phonological conditioning may be contrasted with the principle that determines the selection of yet another allomorph of the past participle morpheme. The final /n/ of *show-n* or *see-n* (which marks them as past participles) is not determined by the phonological structure of the morphs *show* and *see*. For

each English word that is similar to "show" and "see" in this respect, it must be stated as a synchronically inexplicable fact that it selects the /n/ allomorph. This is called grammatical conditioning. There are various kinds of grammatical conditioning.

Alternation of the kind illustrated above for the allomorphs of the plural morpheme and the /ɪd/, /d/, and /t/ allomorphs of the past participle is frequently referred to as morphophonemic. Some linguists have suggested that it should be accounted for not by setting up three allomorphs each with a distinct phonemic form but by setting up a single morph in an intermediate morphophonemic representation. Thus, the regular plural morph might be said to be composed of the morphophoneme /Z/ and the most common past-participle morph of the morphophoneme /D/. General rules of morphophonemic interpretation would then convert /Z/ and /D/ to their appropriate phonetic form according to context. This treatment of the question foreshadows, on the one hand, the stratificational treatment and, on the other, the generative approach, though they differ considerably in other respects.

An important concept in grammar and, more particularly, in morphology is that of free and bound forms. A bound form is one that cannot occur alone as a complete utterance (in some normal context of use). For example, *-ing* is bound in this sense, whereas *wait* is not, nor is *waiting*. Any form that is not bound is free. Bloomfield based his definition of the word on this distinction between bound and free forms. Any free form consisting entirely of two or more smaller free forms was said to be a phrase (e.g., "poor John" or "ran away"), and phrases were to be handled within syntax. Any free form that was not a phrase was defined to be a word and to fall within the scope of morphology. One of the consequences of Bloomfield's definition of the word was that morphology became the study of constructions involving bound forms. The so-called isolating languages, which make no use of bound forms (e.g., Vietnamese), would have no morphology.

The principal division within morphology is between inflection and derivation (or word formation). Roughly speaking, inflectional constructions can be defined as yielding sets of forms that are all grammatically distinct forms of single vocabulary items, whereas derivational constructions yield distinct vocabulary items. For example, "sings," "singing," "sang," and "sung" are all inflectional forms of the vocabulary item traditionally referred to as "the verb to sing"; but "singer," which is formed from "sing" by the addition of the morph *-er* (just as "singing" is formed by the addition of *-ing*), is one of the forms of a different vocabulary item. When this rough distinction between derivation and inflection is made more precise, problems occur. The principal consideration, undoubtedly, is that inflection is more closely integrated with and determined by syntax. But the various formal criteria that have been proposed to give effect to this general principle are not uncommonly in conflict in particular instances, and it probably must be admitted that the distinction between derivation and inflection, though clear enough in most cases, is in the last resort somewhat arbitrary.

Bloomfield and most linguists have discussed morphological constructions in terms of processes. Of these, the most widespread throughout the languages of the world is affixation; i.e., the attachment of an affix to a base. For example, the word "singing" can be described as resulting from the affixation of *-ing* to the base *sing*. (If the affix is put in front of the base, it is a prefix; if it is put after the base, it is a suffix; and if it is inserted within the base, splitting it into two discontinuous parts, it is an infix.) Other morphological processes recognized by linguists need not be mentioned here, but reference may be made to the fact that many of Bloomfield's followers from the mid-1940s were dissatisfied with the whole notion of morphological processes. Instead of saying that *-ing* was affixed to *sing* they preferred to say that *sing* and *-ing* co-occurred in a particular pattern or arrangement, thereby avoiding the implication that *sing* is in some sense prior to or more basic than *-ing*. The distinction of morpheme and morph (and the notion of allomorphs) was developed in order to make possible the description of the morphology and

Plural
morphs in
English

Past
participle
morphs in
English

Inflection
and
derivation

syntax of a language in terms of "arrangements" of items rather than in terms of "processes" operating upon more basic items. Nowadays, the opposition to "processes" is, except among the stratificationists, almost extinct. It has proved to be cumbersome, if not impossible, to describe the relationship between certain linguistic forms without deriving one from the other or both from some common underlying form, and most linguists no longer feel that this is in any way reprehensible.

Bloomfield's theory of syntax

Syntax. Syntax, for Bloomfield, was the study of free forms that were composed entirely of free forms. Central to his theory of syntax were the notions of form classes and constituent structure. (These notions were also relevant, though less central, in the theory of morphology.) Bloomfield defined form classes, rather imprecisely, in terms of some common "recognizable phonetic or grammatical feature" shared by all the members. He gave as examples the form class consisting of "personal substantive expressions" in English (defined as "the forms that, when spoken with exclamatory final pitch, are calls for a person's presence or attention"—e.g., "John," "Boy," "Mr. Smith"); the form class consisting of "infinitive expressions" (defined as "forms which, when spoken with exclamatory final pitch, have the meaning of a command"—e.g., "run," "jump," "come here"); the form class of "nominative substantive expressions" (e.g., "John," "the boys"); and so on. It should be clear from these examples that form classes are similar to, though not identical with, the traditional parts of speech and that one and the same form can belong to more than one form class.

What Bloomfield had in mind as the criterion for form class membership (and therefore of syntactic equivalence) may best be expressed in terms of substitutability. Form classes are sets of forms (whether simple or complex, free or bound), any one of which may be substituted for any other in a given construction or set of constructions throughout the sentences of the language.

The smaller forms into which a larger form may be analyzed are its constituents, and the larger form is a construction. For example, the phrase "poor John" is a construction analyzable into, or composed of, the constituents "poor" and "John." Because there is no intermediate unit of which "poor" and "John" are constituents that is itself a constituent of the construction "poor John," the forms "poor" and "John" may be described not only as constituents but also as immediate constituents of "poor John." Similarly, the phrase "lost his watch" is composed of three word forms—"lost," "his," and "watch"—all of which may be described as constituents of the construction. Not all of them, however, are its immediate constituents. The forms "his" and "watch" combine to make the intermediate construction "his watch"; it is this intermediate unit that combines with "lost" to form the larger phrase "lost his watch." The immediate constituents of "lost his watch" are "lost" and "his watch"; the immediate constituents of "his watch" are the forms "his" and "watch." By the constituent structure of a phrase or sentence is meant the hierarchical organization of the smallest forms of which it is composed (its ultimate constituents) into layers of successively more inclusive units. Viewed in this way, the sentence "Poor John lost his watch" is more than simply a sequence of five word forms associated with a particular intonation pattern. It is analyzable into the immediate constituents "poor John" and "lost his watch," and each of these phrases is analyzable into its own immediate constituents and so on, until, at the last stage of the analysis, the ultimate constituents of the sentence are reached. The constituent structure of the whole sentence is represented by means of a tree diagram in Figure 1.

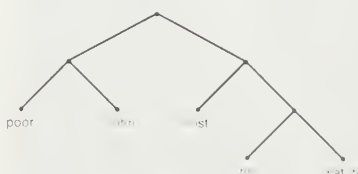


Figure 1: The constituent structure of a simple sentence (see text).

Constituents and immediate constituents

Each form, whether it is simple or composite, belongs to a certain form class. Using arbitrarily selected letters to denote the form classes of English, "poor" may be a member of the form class *A*, "John" of the class *B*, "lost" of the class *C*, "his" of the class *D*, and "watch" of the class *E*. Because "poor John" is syntactically equivalent to (i.e., substitutable for) "John," it is to be classified as a member of *A*. So too, it can be assumed, is "his watch." In the case of "lost his watch" there is a problem. There are very many forms—including "lost," "ate," and "stole"—that can occur, as here, in constructions with a member of *B* and can also occur alone; for example, "lost" is substitutable for "stole the money," as "stole" is substitutable for either or for "lost his watch." This being so, one might decide to classify constructions like "lost his watch" as members of *C*. On the other hand, there are forms that—though they are substitutable for "lost," "ate," "stole," and so on when these forms occur alone—cannot be used in combination with a following member of *B* (cf. "died," "existed"); and there are forms that, though they may be used in combination with a following member of *B*, cannot occur alone (cf. "enjoyed"). The question is whether one respects the traditional distinction between transitive and intransitive verb forms. It may be decided, then, that "lost," "stole," "ate" and so forth belong to one class, *C* (the class to which "enjoyed" belongs), when they occur "transitively" (i.e., with a following member of *B* as their object) but to a different class, *F* (the class to which "died" belongs), when they occur "intransitively." Finally, it can be said that the whole sentence "Poor John lost his watch" is a member of the form class *G*. Thus the constituent structure not only of "Poor John lost his watch" but of a whole set of English sentences can be represented by means of the tree diagram given in Figure 2. New sentences of the same type can be constructed by substituting actual forms for the class labels.

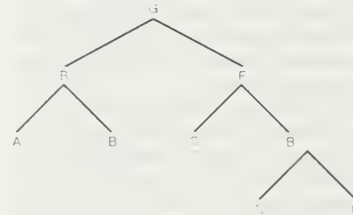


Figure 2: The constituent structure of a class of simple sentences with arbitrary letters used to represent the form class of each constituent (see text).

Any construction that belongs to the same form class as at least one of its immediate constituents is described as endocentric; the only endocentric construction in the model sentence above is "poor John." All the other constructions, according to the analysis, are exocentric. This is clear from the fact that in Figure 2 the letters at the nodes above every phrase other than the phrase *A + B* (i.e., "poor John," "old Harry," and so on) are different from any of the letters at the ends of the lower branches connected directly to these nodes. For example, the phrase *D + E* (i.e., "his watch," "the money," and so forth) has immediately above it a node labelled *B*, rather than either *D* or *E*. Endocentric constructions fall into two types: subordinating and coordinating. If attention is confined, for simplicity, to constructions composed of no more than two immediate constituents, it can be said that subordinating constructions are those in which only one immediate constituent is of the same form class as the whole construction, whereas coordinating constructions are those in which both constituents are of the same form class as the whole construction. In a subordinating construction (e.g., "poor John"), the constituent that is syntactically equivalent to the whole construction is described as the head, and its partner is described as the modifier; thus, in "poor John," the form "John" is the head, and "poor" is its modifier. An example of a coordinating construction is "men and women," in which, it may be assumed, the immediate constituents are the word "men" and the word

Endocentric and exocentric constructions

"women," each of which is syntactically equivalent to "men and women." (It is here implied that the conjunction "and" is not a constituent, properly so called, but an element that, like the relative order of the constituents, indicates the nature of the construction involved. Not all linguists have held this view.)

Ambiguous
construc-
tions

One reason for giving theoretical recognition to the notion of constituent is that it helps to account for the ambiguity of certain constructions. A classic example is the phrase "old men and women," which may be interpreted in two different ways according to whether one associates "old" with "men and women" or just with "men." Under the first of the two interpretations, the immediate constituents are "old" and "men and women"; under the second, they are "old men" and "women." The difference in meaning cannot be attributed to any one of the ultimate constituents but results from a difference in the way in which they are associated with one another. Ambiguity of this kind is referred to as syntactic ambiguity. Not all syntactic ambiguity is satisfactorily accounted for in terms of constituent structure.

Semantics. Bloomfield thought that semantics, or the study of meaning, was the weak point in the scientific investigation of language and would necessarily remain so until the other sciences whose task it was to describe the universe and man's place in it had advanced beyond their present state. In his textbook *Language* (1933), he had himself adopted a behaviouristic theory of meaning, defining the meaning of a linguistic form as "the situation in which the speaker utters it and the response which it calls forth in the hearer." Furthermore, he subscribed, in principle at least, to a physicalist thesis, according to which all science should be modelled upon the so-called exact sciences and all scientific knowledge should be reducible, ultimately, to statements made about the properties of the physical world. The reason for his pessimism concerning the prospects for the study of meaning was his feeling that it would be a long time before a complete scientific description of the situations in which utterances were produced and the responses they called forth in their hearers would be available. At the time that Bloomfield was writing, physicalism was more widely held than it is today, and it was perhaps reasonable for him to believe that linguistics should eschew mentalism and concentrate upon the directly observable. As a result, for some 30 years after the publication of Bloomfield's textbook, the study of meaning was almost wholly neglected by his followers; most American linguists who received their training during this period had no knowledge of, still less any interest in, the work being done elsewhere in semantics.

Result of
Bloom-
field's
physicalist
tendencies

Two groups of scholars may be seen to have constituted an exception to this generalization: anthropologically minded linguists and linguists concerned with Bible translation. Much of the description of the indigenous languages of America has been carried out since the days of Boas and his most notable pupil Sapir by scholars who were equally proficient both in anthropology and in descriptive linguistics; such scholars have frequently added to their grammatical analyses of languages some discussion of the meaning of the grammatical categories and of the correlations between the structure of the vocabularies and the cultures in which the languages operated. It has already been pointed out that Boas and Sapir and, following them, Whorf were attracted by Humboldt's view of the interdependence of language and culture and of language and thought. This view was quite widely held by American anthropological linguists (although many of them would not go as far as Whorf in asserting the dependence of thought and conceptualization upon language).

Also of considerable importance in the description of the indigenous languages of America has been the work of linguists trained by the American Bible Society and the Summer Institute of Linguistics, a group of Protestant missionary linguists. Because their principal aim is to produce translations of the Bible, they have necessarily been concerned with meaning as well as with grammar and phonology. This has tempered the otherwise fairly orthodox Bloomfieldian approach characteristic of the group.

The two most important developments evident in recent

work in semantics are, first, the application of the structural approach to the study of meaning and, second, a better appreciation of the relationship between grammar and semantics. The second of these developments will be treated in the following section on *Transformational-generative grammar*. (See also the separate section below, *Semantics*.) The first, structural semantics, goes back to the period preceding World War II and is exemplified in a large number of publications, mainly by German scholars—Jost Trier, Leo Weisgerber, and their collaborators.

The structural approach to semantics is best explained by contrasting it with the more traditional "atomistic" approach, according to which the meaning of each word in the language is described, in principle, independently of the meaning of all other words. The structuralist takes the view that the meaning of a word is a function of the relationships it contracts with other words in a particular lexical field, or subsystem, and that it cannot be adequately described except in terms of these relationships. For example, the colour terms in particular languages constitute a lexical field, and the meaning of each term depends upon the place it occupies in the field. Although the denotation of each of the words "green," "blue," and "yellow" in English is somewhat imprecise at the boundaries, the position that each of them occupies relative to the other terms in the system is fixed: "green" is between "blue" and "yellow," so that the phrases "greenish yellow" or "yellowish green" and "bluish green" or "greenish blue" are used to refer to the boundary areas. Knowing the meaning of the word "green" implies knowing what cannot as well as what can be properly described as green (and knowing of the borderline cases that they are borderline cases). Languages differ considerably as to the number of basic colour terms that they recognize, and they draw boundaries within the psychophysical continuum of colour at different places. Blue, green, yellow, and so on do not exist as distinct colours in nature, waiting to be labelled differently, as it were, by different languages; they come into existence, for the speakers of particular languages, by virtue of the fact that those languages impose structure upon the continuum of colour and assign to three of the areas thus recognized the words "blue," "green," "yellow."

Colour
terms in
different
languages

The language of any society is an integral part of the culture of that society, and the meanings recognized within the vocabulary of the language are learned by the child as part of the process of acquiring the culture of the society in which he is brought up. Many of the structural differences found in the vocabularies of different languages are to be accounted for in terms of cultural differences. This is especially clear in the vocabulary of kinship (to which a considerable amount of attention has been given by anthropologists and linguists), but it holds true of many other semantic fields also. A consequence of the structural differences that exist between the vocabularies of different languages is that, in many instances, it is in principle impossible to translate a sentence "literally" from one language to another.

It is important, nevertheless, not to overemphasize the semantic incommensurability of languages. Presumably, there are many physiological and psychological constraints that, in part at least, determine one's perception and categorization of the world. It may be assumed that, when one is learning the denotation of the more basic words in the vocabulary of one's native language, attention is drawn first to what might be called the naturally salient features of the environment and that one is, to this degree at least, predisposed to identify and group objects in one way rather than another. It may also be that human beings are genetically endowed with rather more specific and linguistically relevant principles of categorization. It is possible that, although languages differ in the number of basic colour categories that they distinguish, there is a limited number of hierarchically ordered basic colour categories from which each language makes its selection and that what counts as a typical instance, or focus, of these universal colour categories is fixed and does not vary from one language to another. If this hypothesis is correct, then it is false to say, as many structural semanticists have said, that languages divide the continuum of colour in a quite

arbitrary manner. But the general thesis of structuralism is unaffected, for it still remains true that each language has its own unique semantic structure even though the total structure is, in each case, built upon a substructure of universal distinctions.

TRANSFORMATIONAL-GENERATIVE GRAMMAR

A generative grammar, in the sense in which Noam Chomsky uses the term, is a rules system formalized with mathematical precision that generates, without need of any information that is not represented explicitly in the system, the grammatical sentences of the language that it describes, or characterizes, or assigns to each sentence a structural description, or grammatical analysis. All the concepts introduced in this definition of "generative" grammar will be explained and exemplified in the course of this section. Generative grammars fall into several types; this exposition is concerned mainly with the type known as transformational (or, more fully, transformational-generative). Transformational grammar was initiated by Zellig S. Harris in the course of work on what he called discourse analysis (the formal analysis of the structure of continuous text). It was further developed and given a somewhat different theoretical basis by Chomsky.

Harris's grammar. Harris distinguished within the total set of grammatical sentences in a particular language (for example, English) two complementary subsets: kernel sentences (the set of kernel sentences being described as the kernel of the grammar) and nonkernel sentences. The difference between these two subsets lies in nonkernel sentences being derived from kernel sentences by means of transformational rules. For example, "The workers rejected the ultimatum" is a kernel sentence that may be transformed into the nonkernel sentences "The ultimatum was rejected by the workers" or "Did the workers reject the ultimatum?" Each of these may be described as a transform of the kernel sentence from which it is derived. The transformational relationship between corresponding active and passive sentences (e.g., "The workers rejected the ultimatum" and "The ultimatum was rejected by the workers") is conventionally symbolized by the rule $N_1 V N_2 \rightarrow N_2 be V + en$ by N_1 , in which N stands for any noun or noun phrase, V for any transitive verb, en for the past participle morpheme, and the arrow (\rightarrow) instructs one to rewrite the construction to its left as the construction to the right. (There has been some simplification of the rule as it was formulated by Harris.) This rule may be taken as typical of the whole class of transformational rules in Harris's system: it rearranges constituents (what was the first nominal, or noun, N_1 , in the kernel sentence is moved to the end of the transform, and what was the second nominal, N_2 , in the kernel sentence is moved to initial position in the transform), and it adds various elements in specified positions (be, en, and by). Other operations carried out by transformational rules include the deletion of constituents; e.g., the entire phrase "by the workers" is removed from the sentence "The ultimatum was rejected by the workers" by a rule symbolized as $N_2 be V+en$ by $N_1 \rightarrow N_2 be V+en$. This transforms the construction on the left side of the arrow (which resulted from the passive transformation) by dropping the by-phrase, thus producing "The ultimatum was rejected."

Chomsky's grammar. Chomsky's system of transformational grammar, though it was developed on the basis of his work with Harris, differs from Harris's in a number of respects. It is Chomsky's system that has attracted the most attention and has received the most extensive exemplification and further development. As outlined in *Syntactic Structures* (1957), it comprised three sections, or components: the phrase-structure component, the transformational component, and the morphophonemic component. Each of these components consisted of a set of rules operating upon a certain "input" to yield a certain "output." The notion of phrase structure may be dealt with independently of its incorporation in the larger system. In the following system of rules, S stands for Sentence, NP for Noun Phrase, VP for Verb Phrase, Det for Determiner, Aux for Auxiliary (verb), N for Noun, and V for Verb stem.

- (1) S \rightarrow NP + VP
- (2) VP \rightarrow Verb + NP
- (3) NP \rightarrow Det + N
- (4) Verb \rightarrow Aux + V
- (5) Det \rightarrow *the, a, ...*
- (6) N \rightarrow *man, ball, ...*
- (7) Aux \rightarrow *will, can, ...*
- (8) V \rightarrow *hit, see, ...*

This is a simple phrase-structure grammar. It generates and thereby defines as grammatical such sentences as "The man will hit the ball," and it assigns to each sentence that it generates a structural description. The kind of structural description assigned by a phrase-structure grammar is, in fact, a constituent structure analysis of the sentence.

In these rules, the arrow can be interpreted as an instruction to rewrite (this is to be taken as a technical term) whatever symbol appears to the left of the arrow as the symbol or string of symbols that appears to the right of the arrow. For example, rule (2) rewrites the symbol VP as the string of symbols Verb + NP, and it thereby defines Verb + NP to be a construction of the type VP. Or, alternatively and equivalently, it says that constructions of the type VP may have as their immediate constituents constructions of the type Verb and NP (combined in that order). Rule (2) can be thought of as creating or being associated with the tree structure in Figure 3.

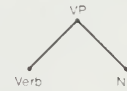


Figure 3: The constituent structure, or phrase structure, assigned by the rule $VP \rightarrow Verb + NP$ (see text).

Rules (1)–(8) do not operate in isolation but constitute an integrated system. The symbol S (standing mnemonically for "sentence") is designated as the initial symbol. This information is not given in the rules (1)–(8), but it can be assumed either that it is given in a kind of protocol statement preceding the grammatical rules or that there is a universal convention according to which S is always the initial symbol. It is necessary to begin with a rule that has the initial symbol on the left. Thereafter any rule may be applied in any order until no further rule is applicable; in doing so, a derivation can be constructed of one of the sentences generated by the grammar. If the rules are applied in the following order: (1), (2), (3), (3), (4), (5), (5), (6), (6), (7), (8), then assuming that "the" is selected on both applications of (5), "man" on one application of (6), and "ball" on the other, "will" on the application of (7), and "hit" on the application of (8), the following derivation of the sentence "The man will hit the ball" will have been constructed:

- | | | |
|--------|---|-------------|
| (i) | S | |
| (ii) | NP + VP | by rule (1) |
| (iii) | NP + Verb + NP | by rule (2) |
| (iv) | Det + N + Verb + NP | by rule (3) |
| (v) | Det + N + Verb + Det + N | by rule (3) |
| (vi) | Det + N + Aux + V + Det + N | by rule (4) |
| (vii) | <i>the</i> + N + Aux + V + Det + N | by rule (5) |
| (viii) | <i>the</i> + N + Aux + V + <i>the</i> + N | by rule (5) |
| (ix) | <i>the</i> + <i>man</i> + Aux + V + <i>the</i> + N | by rule (6) |
| (x) | <i>the</i> + <i>man</i> + Aux + V + <i>the</i> + <i>ball</i> | by rule (6) |
| (xi) | <i>the</i> + <i>man</i> + <i>will</i> + V + <i>the</i> + <i>ball</i> | by rule (7) |
| (xii) | <i>the</i> + <i>man</i> + <i>will</i> + <i>hit</i> + <i>the</i> + <i>ball</i> | by rule (8) |

Many other derivations of this sentence are possible, depending on the order in which the rules are applied. The important point is that all these different derivations are equivalent in that they can be reduced to the same tree diagram: namely, the one shown in Figure 4. If this is compared with the system of rules, it will be seen that each application of each rule creates or is associated with a portion (or subtree) of the tree. The tree diagram, or phrase marker, may now be considered as a structural description of the sentence "The man hit the ball." It is a description of the constituent structure, or phrase structure, of the sentence, and it is assigned by the rules that generate the sentence.

It is important to interpret the term generate in a static, rather than a dynamic, sense. The statement that the

Generation of grammatical sentences

Functions of Harris's rules

Integrated system of rules

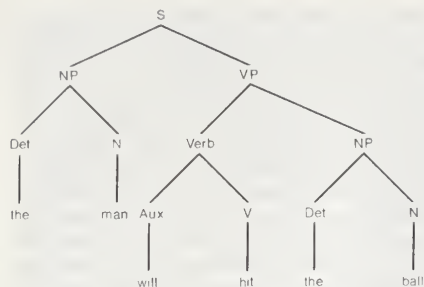


Figure 4: Structural description of the sentence "The man will hit the ball," assigned by the rules of a simple phrase-structure grammar (see text).

grammar generates a particular sentence means that the sentence is one of the totality of sentences that the grammar defines to be grammatical or well formed. All the sentences are generated, as it were, simultaneously. The notion of generation must be interpreted as would be a mathematical formula containing variables. For example, in evaluating the formula $y^2 + y$ for different values of y , one does not say that the formula itself generates these various resultant values (2, when $y = 1$; 5, when $y = 2$; etc.) one after another or at different times; one says that the formula generates them all simultaneously or, better still perhaps, timelessly. The situation is similar for a generative grammar. Although one sentence rather than another can be derived on some particular occasion by making one choice rather than another at particular places in the grammar, the grammar must be thought of as generating all sentences statically or timelessly.

It has been noted that, whereas a phrase-structure grammar is one that consists entirely of phrase-structure rules, a transformational grammar (as formalized by Chomsky) includes both phrase-structure and transformational rules (as well as morphophonemic rules). The transformational rules depend upon the prior application of the phrase-structure rules and have the effect of converting, or transforming, one phrase marker into another. What is meant by this statement may be clarified first with reference to a purely abstract and very simple transformational grammar, in which the letters stand for constituents of a sentence (and S stands for "sentence"):

PS rules

- (1) $S \rightarrow A + B$
- (2) $B \rightarrow C + D$
- (3) $A \rightarrow a + b$
- (4) $C \rightarrow c + e + f$
- (5) $D \rightarrow d + g + h$

T rules

- (6) $A + C + D \rightarrow D + A$

The first five rules are phrase-structure rules (PS rules); rule (6) is a transformational rule (T rule). The output of rules (1)–(5) is the terminal string $a + b + c + e + f + d + g + h$, which has associated with it the structural description indicated by the phrase marker shown in Figure 5 (left). Rule (6) applies to this terminal string of the PS rules and the associated phrase marker. It has the effect of deleting C (and the constituents of C) and permuting A and D (together with their constituents). The result is the string of symbols $d + g + h + a + b$, with the associated phrase marker shown in Figure 5 (right).

The phrase marker shown in Figure 5 (left) may be described as underlying, and the phrase marker shown in

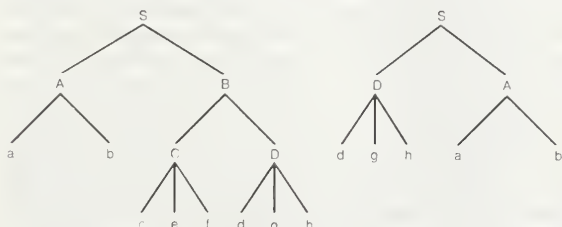


Figure 5: Phrase markers. (Left) A phrase marker associated with the terminal string of a set of phrase-structure rules. (Right) The phrase marker resulting from the application to the phrase marker shown at left of the transformational rule $A + C + D \rightarrow D + A$ (see text).

Figure 5 (right) as derived with respect to rule (6). One of the principal characteristics of a transformational rule is its transformation of an underlying phrase marker into a derived phrase marker in this way. Transformational rules, in contrast with phrase-structure rules, are also formally more heterogeneous and may have more than one symbol on the left-hand side of the arrow. The linguistic importance of these abstract considerations may be explained with reference to the relationship that holds in English between active and passive sentences.

Chomsky's rule for relating active and passive sentences (as given in *Syntactic Structures*) is very similar, at first sight, to Harris's, discussed above. Chomsky's rule is:

$$NP_1 - Aux - V - NP_2 \rightarrow NP_2 - Aux + be + en - V - by + NP_1$$

This rule, called the passive transformation, presupposes and depends upon the prior application of a set of phrase-structure rules. For simplicity, the passive transformation may first be considered in relation to the set of terminal strings generated by the phrase-structure rules (1)–(8) given earlier. The string "the + man + will + hit + the + ball" (with its associated phrase marker, as shown in Figure 4) can be treated not as an actual sentence but as the structure underlying both the active sentence "The man will hit the ball" and the corresponding passive "The ball will be hit by the man." The passive transformation is applicable under the condition that the underlying, or "input," string is analyzable in terms of its phrase structure as $NP - Aux - V - NP$ (the use of subscript numerals to distinguish the two NPs in the formulation of the rule is an informal device for indicating the operation of permutation). In the phrase marker in Figure 4 "the" + "man" are constituents of NP, "will" is a constituent of Aux, "hit" is a constituent of V, and "the" + "ball" are constituents of NP. The whole string is therefore analyzable in the appropriate sense, and the passive transformation converts it into the string "the + ball + will + be + en + hit + by + the + man." A subsequent transformational rule will permute "en + hit" to yield "hit + en," and one of the morphophonemic rules will then convert "hit + en" to "hit" (as "ride + en" will be converted to "ridden"; "open + en" to "opened," and so on).

Passive transformation

Every transformational rule has the effect of converting an underlying phrase marker into a derived phrase marker. The manner in which the transformational rules assign derived constituent structure to their input strings is one of the major theoretical problems in the formalization of transformational grammar. Here it can be assumed not only that "be + en" is attached to Aux and "by" to NP (as indicated by the plus signs in the rule as it has been formulated above) but also that the rest of the derived structure is as shown in Figure 6. The phrase marker in Figure 6 formalizes the fact, among others, that "the ball" is the subject of the passive sentence "The ball will be hit by the man," whereas "the man" is the subject of the corresponding active "The man will hit the ball" (cf. Figure 4).

Although the example above is a very simple one, and only a single transformational rule has been considered independently of other transformational rules in the same system, the passive transformation must operate, not only upon simple noun phrases like "the man" or "the ball," but upon noun phrases that contain adjectives ("the old man"), modifying phrases ("the man in the corner"), relative clauses ("the man who checked in last night"), and so forth. The incorporation, or embedding, of these other structures with the noun phrase will be brought about by the prior application of other transformational rules. It should also be clear that the phrase-structure rules require extension to allow for the various forms of the verb ("is hitting," "hit," "was hitting," "has hit," "has been hitting," etc.) and for the distinction of singular and plural.

It is important to note that, unlike Harris's, Chomsky's system of transformational grammar does not convert one sentence into another: the transformational rules operate upon the structures underlying sentences and not upon actual sentences. A further point is that even the simplest sentences (*i.e.*, kernel sentences) require the application of

Operation of transformations on underlying structures

Phrase-structure and transformational rules

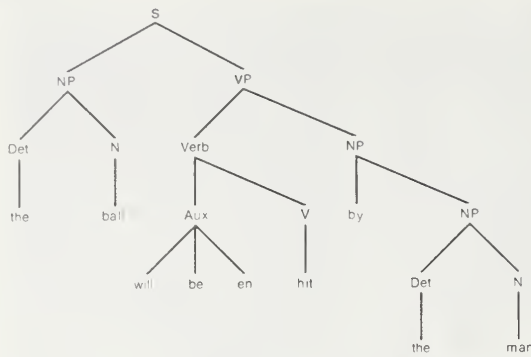


Figure 6: A possible derived phrase marker for a passive sentence (see text).

at least some transformational rules. Corresponding active and passive sentences, affirmative and negative sentences, declarative and interrogative sentences, and so on are formally related by deriving them from the same underlying terminal string of the phrase-structure component. The difference between kernel sentences and nonkernel sentences in *Syntactic Structures* (in Chomsky's later system the category of kernel sentences is not given formal recognition at all) resides in the fact that kernel sentences are generated without the application of any optional transformations. Nonkernel sentences require the application of both optional and obligatory transformations, and they differ one from another in that a different selection of optional transformations is made.

Modifications in Chomsky's grammar. Chomsky's system of transformational grammar was substantially modified in 1965. Perhaps the most important modification was the incorporation, within the system, of a semantic component, in addition to the syntactic component and phonological component. (The phonological component may be thought of as replacing the morphophonemic component of *Syntactic Structures*.) The rules of the syntactic component generate the sentences of the language and assign to each not one but two structural analyses: a deep structure analysis as represented by the underlying phrase marker, and a surface structure analysis, as represented by the final derived phrase marker. The underlying phrase marker is assigned by rules of the base (roughly equivalent to the PS [Phrase-Structure] rules of the earlier system); the derived phrase marker is assigned by the transformational rules. The interrelationship of the four sets of rules is shown diagrammatically in Figure 7. The

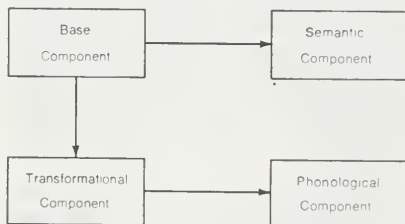


Figure 7: Diagrammatic representation of a transformational grammar (see text).

meaning of the sentence is derived (mainly, if not wholly) from the deep structure by means of the rules of semantic interpretation; the phonetic realization of the sentence is derived from its surface structure by means of the rules of the phonological component. The grammar ("grammar" is now to be understood as covering semantics and phonology, as well as syntax) is thus an integrated system of rules for relating the pronunciation of a sentence to its meaning. The syntax, and more particularly the base, is at the "heart" of the system, as it were: it is the base component (as the arrows in the diagram indicate) that generates the infinite class of structures underlying the well-formed sentences of a language. These structures are then given a semantic and phonetic "interpretation" by the other components.

The base consists of two parts: a set of categorial rules

and a lexicon. Taken together, they fulfill a similar function to that fulfilled by the phrase-structure rules of the earlier system. But there are many differences of detail. Among the most important is that the lexicon (which may be thought of as a dictionary of the language cast in a particular form) lists, in principle, all the vocabulary words in the language and associates with each all the syntactic, semantic, and phonological information required for the correct operation of the rules. This information is represented in terms of what are called features. For example, the entry for "boy" might say that it has the syntactic features: [+ Noun], [+ Count], [+ Common], [+ Animate], and [+ Human]. The categorial rules generate a set of phrase markers that have in them, as it were, a number of "slots" to be filled with items from the lexicon. With each such "slot" there is associated a set of features that define the kind of item that can fill the "slot." If a phrase marker is generated with a "slot" for the head of a noun phrase specified as requiring an animate noun (*i.e.*, a noun having the feature [+ Animate]), the item "boy" would be recognized as being compatible with this specification and could be inserted in the "slot" by the rule of lexical substitution. Similarly, it could be inserted in "slots" specified as requiring a common noun, a human noun, or a countable noun, but it would be excluded from positions that require an abstract noun (*e.g.*, "sincerity") or an uncountable noun (*e.g.*, "water"). By drawing upon the syntactic information coded in feature notation in the lexicon, the categorial rules might permit such sentences as "The boy died," while excluding (and thereby defining as ungrammatical) such nonsentences as "The boy elapsed."

One of the most controversial topics in the development of transformational grammar is the relationship between syntax and semantics. Scholars working in the field are now agreed that there is a considerable degree of interdependence between the two, and the problem is how to formalize this interdependence. One school of linguists, called generative semanticists, accept the general principles of transformational grammar but have challenged Chomsky's conception of deep structure as a separate and identifiable level of syntactic representation. In their opinion, the basic component of the grammar should consist of a set of rules for the generation of well-formed semantic representations. These would then be converted by a succession of transformational rules into strings of words with an assigned surface-structure syntactic analysis, there being no place in the passage from semantic representation to surface structure identifiable as Chomsky's deep structure. Chomsky himself has denied that there is any real difference between the two points of view and has maintained that the issue is purely one of notation. That this argument can be put forward by one party to the controversy and rejected by the other is perhaps a sufficient indication of the uncertainty of the evidence. Of greater importance than the overt issues, in so far as they are clear, is the fact that linguists are now studying much more intensively than they have in the past the complexities of the interdependence of syntax, on the one hand, and semantics and logic, on the other. Whether it will prove possible to handle all these complexities within a comprehensive generative grammar remains to be seen.

The role of the phonological component of a generative grammar of the type outlined by Chomsky is to assign a phonetic "interpretation" to the strings of words generated by the syntactic component. These strings of words are represented in a phonological notation (taken from the lexicon) and have been provided with a surface-structure analysis by the transformational rules (see Figure 7). The phonological elements out of which the word forms are composed are segments consisting of what are referred to technically as distinctive features (following the usage of the Prague school, see below). For example, the word form "man," represented phonologically, is composed of three segments: the first consists of the features [+ consonantal], [+ bilabial], [+ nasal], etc.; the second of the features [+ vocalic], [+ front], [+ open], etc.; and the third of the features [+ consonantal], [+ alveolar], [+ nasal], etc. (These features should be taken as purely illustrative; there is some doubt about the definitive list of distinctive

Base:
categorial
rules and a
lexicon

Generative
semantic-
tists

features.) Although these segments may be referred to as the "phonemes" /m/, /a/, and /n/, they should not be identified theoretically with units of the kind discussed in the section on *Phonology* under *Structural linguistics*. They are closer to what many American structural linguists called "morphophonemes" or the Prague school linguists labelled "archiphonemes," being unspecified for any feature that is contextually redundant or predictable. For instance, the first segment of the phonological representation of "man" will not include the feature [+ voice]; because nasal consonants are always phonetically voiced in this position in English, the feature [+ voice] can be added to the phonetic specification by a rule of the phonological component.

One further important aspect of generative phonology (*i.e.*, phonology carried out within the framework of an integrated generative grammar) should be mentioned: its dependence upon syntax. Most American structural phonologists made it a point of principle that the phonemic analysis of an utterance should be carried out without regard to its grammatical structure. This principle was controversial among American linguists and was not generally accepted outside America. Not only has the principle been rejected by the generative grammarians, but they have made the phonological description of a language much more dependent upon its syntactic analysis than has any other school of linguists. They have claimed, for example, that the phonological rules that assign different degrees of stress to the vowels in English words and phrases and alter the quality of the relatively unstressed vowel concomitantly must make reference to the derived constituent structure of sentences and not merely to the form class of the individual words or the places in which the word boundaries occur.

TAGMEMICS

The system of tagmemic analysis, as presented by Kenneth L. Pike, was developed for the analysis not only of language but of all of human behaviour that manifests the property of patterning. In the following treatment, only language will be discussed.

Modes of language. Every language is said to be trimodal—*i.e.*, structured in three modes: phonology, grammar, and lexicon. These modes are interrelated but have a considerable degree of independence and must be described in their own terms. Phonology and lexicon should not be seen as mere appendages to grammar, the former simply specifying which phonemes can combine to form morphemes (or morphs), and the latter simply listing the morphemes and other meaningful units with a description of their meaning. There are levels of structure in each of the modes, and the units of one level are not necessarily coterminous with those of another. Phonemes, for example, may combine to form syllables and syllables to form phonological words ("phonological word" is defined as the domain of some phonological process such as accentuation, assimilation, or dissimilation), but the morpheme (or morph) will not necessarily consist of an integral number of syllables, still less of a single syllable. Nor will the word as a grammatical unit necessarily coincide with the phonological word. Similarly, the units of lexical analysis, sometimes referred to as lexemes (in one sense of this term), are not necessarily identifiable as single grammatical units, whether as morphemes, words, or phrases. No priority, then, is ascribed to any one of the three modes.

The originality of tagmemic analysis and the application of the term tagmeme is most clearly manifest in the domain of grammar. By a tagmeme is meant an element of a construction, the element in question being regarded as a composite unit, described in such terms as "slot-filler" or "function-class." For example, one of the tagmemes required for the analysis of English at the syntactic level might be noun-as-subject, in which "noun" refers to a class of substitutable, or paradigmatically related, morphemes or words capable of fulfilling a certain grammatical function, and "subject" refers to the function that may be fulfilled by one or more classes of elements. In the tagmeme noun-as-subject—which, using the customary tagmemic symbolism, may be represented as Subject:noun—the sub-

ject slot is filled by a noun. When a particular tagmeme is identified in the analysis of an actual utterance, it is said to be manifested by the particular member of the grammatical class that occurs in the appropriate slot in the utterance. For example, in the utterance "John is asleep," the subject tagmeme is manifested by the noun "John." Tagmemicists insist that tagmemes, despite their bipartite structure, are single units. In grammatical analysis, the distribution of tagmemes, not simply of classes, is stated throughout the sentences of the language. Subject:noun is a different tagmeme from Object:noun, as it is also a different tagmeme from Subject:pronoun.

Hierarchy of levels. Within the grammar of a language there is a hierarchy of levels, units of one level being composed of sequences of units of the level below. In many languages, five such levels are recognized, defined in terms of the following units: morpheme, word, phrase, clause, and sentence. (The term level is being used in a different sense from that in which it was used earlier to refer to phonology and grammar.) The difference between morphology and syntax is simply a difference between two of these five levels, no greater than the difference, for example, between the phrase level and the clause level. Normally, tagmemes at one level are manifested by units belonging to the level below: clause tagmemes by phrases, phrase tagmemes by words, and so on. Intermediate levels may, however, be skipped. For example, the subject tagmeme in a clause may be manifested by a single word in English (*e.g.*, "John," "water") and not necessarily by a phrase ("the young man").

It is also possible for there to be loop-backs in the grammatical hierarchy of a language. This means that a unit of higher level may be embedded within the structure of a unit of lower level; for example, a clause may fill a slot within a phrase (*e.g.*, "who arrived late," in "the man who arrived late").

In regard to the notation of tagmemics, a construction is symbolized as a string of tagmemes (which commonly, though not necessarily, will be sequentially ordered according to the order in which elements manifesting the tagmemes occur in utterances). Each tagmeme is marked as obligatory or optional by having preposed to it a plus sign (+) or a plus-or-minus sign (\pm), respectively. For example, a formula representing the structure of a clause composed of an obligatory subject tagmeme, an obligatory predicate tagmeme, and an optional object tagmeme might be $Cl = + S:n + P:v \pm O:n$ (in which Cl stands for a clause of a certain type and n and v stand for the classes of nouns and verbs, respectively). This formula does not represent in any way the fact (if it is a fact) that the predicate tagmeme and object tagmeme together form a unit that is one of the two immediate constituents of the clause. It is one of the characteristic features of tagmemic grammar that it gives much less emphasis to the notion of constituent structure than other American approaches to grammatical analysis.

STRATIFICATIONAL GRAMMAR

This system of analysis (whose principal advocate is Sydney M. Lamb, a U.S. linguist) is called stratificational because it is based upon the notion that every language comprises a restricted number of structural layers or strata, hierarchically related in such a way that units or combinations of units on one stratum realize units or combinations of units of the next higher stratum. The number of strata may vary from language to language. Four strata have been recognized for English, and it is probable that all languages may have at least these four: the sememic, the lexemic, the morphemic, and the phonemic strata. The sememic stratal system constitutes the semology of the language; the lexemic and morphemic stratal systems constitute the grammar (in the narrower sense of this term); and the phonemic system constitutes the phonology. In some later stratificational work, the term grammar covers the three higher stratal systems—the sememic, the lexemic, and the morphemic—and is opposed to "phonology." The deep structure of sentences is described on the sememic stratum and the surface structure on the morphemic. In the present account, "grammar" is used in

Depen-
dence of
generative
phonology
on syntax

Tagmemic
notation

Tagmeme
defined

Sememic,
lexemic,
morphemic,
and
phonemic
strata

the narrower sense and will be opposed to "semology" as well as "phonology."

The originality of stratificational grammar does not reside in the recognition of these three major components of a linguistic description. The stratificational approach to linguistic description is distinguished from others in that it relates grammar to semology and phonology by means of the same notion of realization that it employs to relate the lexemic and the morphemic stratal systems within the grammatical component. Another distinguishing feature of stratificational grammar, in its later development at least, is its description of linguistic structure in terms of a network of relationships, rather than by means of a system of rules; linguistic units are said to be nothing more than points, or positions, in the relational network.

Technical terminology. Lamb has been very careful to make the terminology of stratificational grammar as consistent and perspicuous as possible; but, in fitting some of the more or less established terms into his own theoretical framework, he has reinterpreted them in a potentially confusing manner. Thus, the same terms have been used in different senses in different versions of the system. For example, "morpheme" in stratificational grammar corresponds neither to the unit to which Bloomfield applied the term (*i.e.*, to a word segment consisting of phonemes) nor to the more abstract grammatical unit that a Bloomfieldian morpheme might be described as representing (*e.g.*, the past-tense morpheme that might be variously represented by such allomorphs as /id/, /t/, /d/, etc.). Lamb describes the morpheme as a unit composed of morphons (roughly equivalent to what other linguists have called morphophonemes) that is related to a combination of one or more compositional units of the stratum above, lexons, by means of the relationship of realization. For example, the word form "hated" realizes (on the morphemic stratum) a combination of two lexons, one of which, the stem, realizes the lexeme HATE and the other, the suffix, realizes the PAST TENSE lexeme; each of these two lexons is realized on the stratum below by a morpheme. Another example brings out more clearly the difference between morphemes (the minimal grammatical elements) and lexemes (the minimal meaningful elements). The word form "understood" realizes a combination of three morphemes UNDER, STAND, and PAST. UNDER and STAND jointly realize the single lexeme UNDERSTAND (whose meaning cannot be described as a function of the meanings of UNDER and STAND), whereas the single PAST morpheme directly realizes the single lexeme PAST TENSE.

The stratificational framework, presented in Lamb's work, consistently separates compositional and realizational units, the former being designated by terms ending in the suffix *-on* (semon, lexon, morphon, phonon), the latter by terms ending in the suffix *-eme* (sememe, lexeme, morpheme, phoneme). Ons are components or compounds of emes on the same stratum (semons are components of sememes, lexons are composed of lexemes, etc.) and emes realize ons of the stratum above (phonemes realize morphons, morphemes realize lexons, etc.). Each stratum has its own combinatorial pattern specifying the characteristic combinations of elements on that stratum. Syllable structure is specified on the phonemic stratum, the structure of word forms on the morphemic stratum, the structure of phrases on the lexemic stratum, and the structure of clauses and sentences on the sememic stratum. Phonons are roughly equivalent to phonological distinctive features and include such properties or components of phonemes as labial, nasal, and so on. Sememes are roughly equivalent to what other linguists have called semantic components or features and include such aspects of the meaning of the lexeme "man" as "male," "adult," "human," and so forth. Once again, however, compositional function is distinguished from interstratal realizational function, so that no direct equivalence can be established with nonstratificational terminology. In more recent work in stratificational grammar, the notion that emes are composed of ons has been abandoned, and greater emphasis is laid upon the fact that emes are points, or positions, in a relational network; they are connected to other points in the network but have themselves no internal structure.

Interstratal relationships. One of the principal characteristics of the stratificational approach is that it sets out to describe languages without making use of rules that convert one entity into another. (Reference has been made above to the antipathy many linguists have felt towards describing languages in terms of processes.) The stratificationalist would handle the phenomena in terms of the interstratal relationships of realization. Various kinds of interstratal relationships, other than that of one-to-one correspondence may be recognized: diversification, in which one higher unit has alternative realizations; zero realization, in which a higher unit has no overt realization on the lower stratum; neutralization, in which two or more higher units are mapped into the same lower level unit; and so on. All these interstratal one-many or many-one relations are then analyzed in terms of the logical notions of conjunction and disjunction (AND-relations versus OR-relations), of ordering (*x* precedes *y* in an AND-relationship, *x* is selected in preference to *y* in an OR-relationship), and the directionality ("upward" towards meaning, or "downward" towards sound). Many of the phenomena that are described by other linguists in terms of processes that derive one unit from another can be described elegantly enough in terms of interstratal relationships of this kind.

Critics, however, have objected to the proliferation of strata and theoretical constructs in stratificational grammar, arguing that they result from an a priori commitment to the notion of realization and that the only stratal distinction for which there is any independent evidence is the distinction of phonology and grammar. It has been suggested by Lamb that stratificational grammar provides a model for the way in which linguistic information is stored in the brain and activated during the production and reception of speech. But little is known yet about the neurology of language and speech, and it would be premature to draw any firm conclusions about this aspect of stratificational grammar (see below *Psycholinguistics*).

THE PRAGUE SCHOOL

What is now generally referred to as the Prague school comprises a fairly large group of scholars, mainly European, who, though they may not themselves have been members of the Linguistic Circle of Prague, derived their inspiration from the work of Vilém Mathesius, Nikolay Trubetsky, Roman Jakobson and other scholars based in Prague in the decade preceding World War II.

Combination of structuralism and functionalism. The most characteristic feature of the Prague school approach is its combination of structuralism with functionalism. The latter term (like "structuralism") has been used in a variety of senses in linguistics. Here it is to be understood as implying an appreciation of the diversity of functions fulfilled by language and a theoretical recognition that the structure of languages is in large part determined by their characteristic functions. Functionalism, taken in this sense, manifests itself in many of the more particular tenets of Prague school doctrine.

One very famous functional analysis of language, which, though it did not originate in Prague, was very influential there, was that of the German psychologist Karl Bühler, who recognized three general kinds of function fulfilled by language: *Darstellungsfunktion*, *Kundgabefunktion*, and *Appellfunktion*. These terms may be translated, in the present context, as the cognitive, the expressive, and the conative (or instrumental) functions. The cognitive function of language refers to its employment for the transmission of factual information; by expressive function is meant the indication of the mood or attitude of the speaker (or writer); and by the conative function of language is meant its use for influencing the person one is addressing or for bringing about some practical effect. A number of scholars working in the Prague tradition have suggested that these three functions correlate in many languages, at least partly, with the grammatical categories of mood and person. The cognitive function is fulfilled characteristically by 3rd-person nonmodal utterances (*i.e.*, utterances in the indicative mood, making no use of modal verbs); the expressive function by 1st-person utterances in the subjunctive or optative mood; and the conative function by 2nd-person

Types of interstratal relationships

Morpheme defined in stratificational terms

Various functions of language

utterances in the imperative. The functional distinction of the cognitive and the expressive aspects of language has also been applied by Prague school linguists in their work on stylistics and literary criticism. One of their key principles is that language is being used poetically or aesthetically when the expressive aspect is predominant, and that it is typical of the expressive function of language that this should be manifest in the form of an utterance and not merely in the meanings of the component words.

Phonological contributions. The Prague school is best known for its work on phonology. Unlike the American phonologists, Trubetsky and his followers did not take the phoneme to be the minimal unit of analysis. Instead, they defined phonemes as sets of distinctive features. For example, in English, /b/ differs from /p/ in the same way that /d/ differs from /t/ and /g/ from /k/. Just how they differ in terms of their articulation is a complex question. For simplicity, it may be said that there is just one feature, the presence of which distinguishes /b/, /d/, and /g/ from /p/, /t/, and /k/, and that this feature is voicing (vibration of the vocal cords). Similarly, the feature of labiality can be extracted from /p/ and /b/ by comparing them with /t/, /d/, /k/, and /g/; the feature of nasality from /n/ and /m/ by comparing them with /t/ and /d/, on the one hand, and with /p/ and /b/, on the other. Each phoneme, then, is composed of a number of articulatory features and is distinguished by the presence or absence of at least one feature from every other phoneme in the language. The distinctive function of phonemes, which depends upon and supports the principle of the duality of structure, can be related to the cognitive function of language. This distinctive feature analysis of Prague school phonology as developed by Jakobson has become part of the generally accepted framework for generative phonology (see above).

Two other kinds of phonologically relevant function are also recognized by linguists of the Prague school: expressive and demarcative. The former term is employed here in the sense in which it was employed above (*i.e.*, in opposition to "cognitive"); it is characteristic of stress, intonation, and other suprasegmental aspects of language that they are frequently expressive of the mood and attitude of the speaker in this sense. The term demarcative is applied to those elements or features that in particular languages serve to indicate the occurrence of the boundaries of words and phrases and, presumably, make it easier to identify such grammatical units in the stream of speech. There are, for example, many languages in which the set of phonemes that can occur at the beginning of a word differs from the set of phonemes that can occur at the end of a word. These and other devices are described by the Prague school phonologists as having demarcative function: they are boundary signals that reinforce the identity and syntagmatic unity of words and phrases.

Theory of markedness. The notion of markedness was first developed in Prague school phonology but was subsequently extended to morphology and syntax. When two phonemes are distinguished by the presence or absence of a single distinctive feature, one of them is said to be marked and the other unmarked for the feature in question. For example, /b/ is marked and /p/ unmarked with respect to voicing. Similarly, in morphology, the regular English verb can be said to be marked for past tense (by the suffixation of *-ed*) but to be unmarked in the present (*cf.* "jumped" versus "jump"). It is often the case that a morphologically unmarked form has a wider range of occurrences and a less definite meaning than a morphologically marked form. It can be argued, for example, that, whereas the past tense form in English (in simple sentences or the main clause of complex sentences) definitely refers to the past, the so-called present tense form is more neutral with respect to temporal reference: it is nonpast in the sense that it fails to mark the time as past, but it does not mark it as present. There is also a more abstract sense of markedness, which is independent of the presence or absence of an overt feature or affix. The words "dog" and "bitch" provide examples of markedness of this kind on the level of vocabulary. Whereas the use of the word "bitch" is restricted to females of the species, "dog" is applicable to both males and females. "Bitch"

is the marked and "dog" the unmarked term, and, as is commonly the case, the unmarked term can be neutral or negative according to context (*cf.* "That dog over there is a bitch" versus "It's not a dog, it's a bitch"). The principle of markedness, understood in this more general or more abstract sense, is now quite widely accepted by linguists of many different schools, and it is applied at all levels of linguistic analysis.

Recent contributions. Current Prague school work is still characteristically functional in the sense in which this term was interpreted in the pre-World War II period. The most valuable contribution made by the postwar Prague school is probably the distinction of theme and rheme and the notion of "functional sentence perspective" or "communicative dynamism." By the theme of a sentence is meant that part that refers to what is already known or given in the context (sometimes called, by other scholars, the topic or psychological subject); by the rheme, the part that conveys new information (the comment or psychological predicate). It has been pointed out that, in languages with a free word order (such as Czech or Latin), the theme tends to precede the rheme, regardless of whether the theme or the rheme is the grammatical subject and that this principle may still operate, in a more limited way, in languages, like English, with a relatively fixed word order (*cf.* "That book I haven't seen before"). But other devices may also be used to distinguish theme and rheme. The rheme may be stressed ("John saw Mary") or made the complement of the verb "to be" in the main clause of what is now commonly called a cleft sentence ("It's John who saw Mary").

The general principle that has guided research in "functional sentence perspective" is that the syntactic structure of a sentence is in part determined by the communicative function of its various constituents and the way in which they relate to the context of utterance. A somewhat different but related aspect of functionalism in syntax is seen in current work in what is called case grammar. Case grammar is based upon a small set of syntactic functions (agentive, locative, benefactive, instrumental, and so on) that are variously expressed in different languages but that are held to determine the grammatical structure of sentences. Although case grammar does not derive directly from the work of the Prague school, it is very similar in inspiration.

Case
grammar

Historical (diachronic) linguistics

LINGUISTIC CHANGE

All languages change in the course of time. Written records make it clear that 15th-century English is quite noticeably different from 20th-century English, as is 15th-century French or German from modern French or German. It was the principal achievement of the 19th-century linguists not only to realize more clearly than their predecessors the ubiquity of linguistic change but also to put its scientific investigation on a sound footing by means of the comparative method (see the section *History of linguistics: The 19th century*). This will be treated in greater detail in the following section. Here various kinds, or categories, of linguistic change will be listed and exemplified.

Sound change. Since the beginning of the 19th century, when scholars observed that there were a number of systematic correspondences in related words between the sounds of the Germanic languages and the sounds of what were later recognized as other Indo-European languages, particular attention has been paid in diachronic linguistics to changes in the sound systems of languages.

Certain common types of sound change, most notably assimilation and dissimilation, can be explained, at least partially, in terms of syntagmatic, or contextual, conditioning. By assimilation is meant the process by which one sound is made similar in its place or manner of articulation to a neighbouring sound. For example, the word "cupboard" was presumably once pronounced, as the spelling indicates, with the consonant cluster *pb* in the middle. The *p* was assimilated to *b* in manner of articulation (*i.e.*, voicing was maintained throughout the cluster), and subsequently the resultant double consonant *bb* was

Phonemes
as
distinctive
features

Marked
and
unmarked
features

simplified. With a single *b* in the middle and an unstressed second syllable, the word "cupboard," as it is pronounced nowadays, is no longer so evidently a compound of "cup" and "board" as its spelling still shows it to have been. The Italian words *notte* "night" and *otto* "eight" manifest assimilation of the first consonant to the second consonant of the cluster in place of articulation (cf. Latin *nocte(m)*, *octo*). Assimilation is also responsible for the phenomenon referred to as umlaut in the Germanic languages. The high front vowel *i* of suffixes had the effect of fronting and raising preceding back vowels and, in particular, of converting an *a* sound into an *e* sound. In Modern German this is still a morphologically productive process (cf. *Mann* "man"; *Männer* "men"). In English it has left its mark in such irregular forms as "men" (from **manniz*), "feet" (from **fotiz*), and "length" (from **langha*).

Umlaut

Dissimilation refers to the process by which one sound becomes different from a neighbouring sound. For example, the word "pilgrim" (French *pèlerin*) derives ultimately from the Latin *peregrinus*; the *l* sound results from dissimilation of the first *r* under the influence of the second *r*. A special case of dissimilation is haplology, in which the second of the two identical or similar syllables is dropped. Examples include the standard modern British pronunciations of "Worcester" and "Gloucester" with two syllables rather than three and the common pronunciation of "library" as if it were written "libry." Both assimilation and dissimilation are commonly subsumed under the principle of "ease of articulation." This is clearly applicable in typical instances of assimilation. It is less obvious how or why a succession of unlike sounds in contiguous syllables should be easier to pronounce than a succession of identical or similar sounds. But a better understanding of this phenomenon, as of other "slips of the tongue," may result from current work in the physiological and neurological aspects of speech production.

Not all sound change is to be accounted for in terms of syntagmatic conditioning. The change of *p*, *t*, and *k* to *f*, *θ* (the *th* sound in "thin"), and *h* or of *b*, *d*, *g* to *p*, *t*, and *k* in early Germanic cannot be explained in these terms. Nor can the so-called Great Vowel Shift in English, which, in the 15th century, modified the quality of all the long vowels (cf. "profane" : "profanity"; "divine" : "divinity"; and others). Attempts have been made to develop a general theory of sound change, notably by the French linguist André Martinet. But no such theory has yet won universal acceptance, and it is likely that the causes of sound change are multiple.

Phonetic vs. phonological sound change

Sound change is not necessarily phonological; it may be merely phonetic (see above *Structural linguistics: Phonology*). The pronunciation of one or more of the phones realizing a particular phoneme may change slightly without affecting any of the previously existing phonological distinctions; this no doubt happens quite frequently as a language is transmitted from one generation to the next. Two diachronically distinct states of the language would differ in this respect in the same way as two coexistent but geographically or socially distinct accents of the same language might differ. It is only when two previously distinct phonemes are merged or a unitary phoneme splits into two (typically when allophonic variation becomes phonemic) that sound change must definitely be considered as phonological. For example, the sound change of *p* to *f*, *t* to *θ* (*th*), and *k* to *h*, on the one hand, and of *b* to *p*, *d* to *t*, and *g* to *k*, on the other, in early Germanic had the effect of changing the phonological system. The voiceless stops did not become fricatives in all positions; they remained as voiceless stops after *s*. Consequently, the *p* sound that was preserved after *s* merged with the *p* that derived by sound change from *b*. (It is here assumed that the aspirated *p* sound and the unaspirated *p* sound are to be regarded as allophones of the same phoneme). Prior to the Germanic sound shift the phoneme to be found at the beginning of the words for "five" or "father" also occurred after *s* in words for "spit" or "spew"; after the change this was no longer the case.

Grammatical change. A language can acquire a grammatical distinction that it did not have previously, as when English developed the progressive ("He is running")

in contrast to the simple present ("He runs"). It can also lose a distinction; e.g., modern spoken French has lost the distinction between the simple past (*Il marcha* "he walked") and the perfect (*Il a marché* "he has walked"). What was expressed by means of one grammatical device may come to be expressed by means of another. For example, in the older Indo-European languages the syntactic function of the nouns and noun phrases in a sentence was expressed primarily by means of case endings (the subject of the sentence being in the nominative case, the object in the accusative case, and so on); in most of the modern Indo-European languages these functions are expressed by means of word order and the use of prepositions. It is arguable, although it can hardly be said to have been satisfactorily demonstrated yet, that the grammatical changes that take place in a language in the course of time generally leave its deep structure unaffected and tend to modify the ways in which the deeper syntactic functions and distinctions are expressed (whether morphologically, by word order, by the use of prepositions and auxiliary verbs, or otherwise), without affecting the functions and distinctions themselves. Many grammatical changes are traditionally accounted for in terms of analogy.

Semantic change. Towards the end of the 19th century, a French scholar, Michel Bréal, set out to determine the laws that govern changes in the meaning of words. This was the task that dominated semantic research until the 1930s, when scholars began to turn their attention to the synchronic study of meaning. Many systems for the classification of changes of meaning have been proposed, and a variety of explanatory principles have been suggested. So far no "laws" of semantic change comparable to the phonologist's sound laws have been discovered. It seems that changes of meaning can be brought about by a variety of causes. Most important, perhaps, and the factor that has been emphasized particularly by the so-called words-and-things movement in historical semantics is the change undergone in the course of time by the objects or institutions that words denote. For example, the English word "car" goes back through Latin *carrus* to a Celtic word for a four-wheeled wagon. It now denotes a very different sort of vehicle; confronted with a model of a Celtic wagon in a museum, one would not describe it as a car.

Causes of changes in meaning

Some changes in the meaning of words are caused by their habitual use in particular contexts. The word "starve" once meant "to die" (cf. Old English *steorfan*, German *sterben*); in most dialects of English, it now has the more restricted meaning "to die of hunger," though in the north of England "He was starving" can also mean "He was very cold" (i.e., "dying" of cold, rather than hunger). Similarly, the word "deer" has acquired a more specialized meaning than the meaning "wild animal" that it once bore (cf. German *Tier*); and "meat," which originally meant food in general (hence, "sweetmeats" and the archaic phrase "meat and drink") now denotes the flesh of an animal treated as food. In all such cases, the narrower meaning has developed from the constant use of the word in a more specialized context, and the contextual presuppositions of the word have in time become part of its meaning.

Borrowing. Languages borrow words freely from one another. Usually this happens when some new object or institution is developed for which the borrowing language has no word of its own. For example, the large number of words denoting financial institutions and operations borrowed from Italian by the other western European languages at the time of the Renaissance testifies to the importance of the Italian bankers in that period. (The word "bank" itself, in this sense, comes through French from the Italian *banca*). Words now pass from one language to another on a scale that is probably unprecedented, partly because of the enormous number of new inventions that have been made in the 20th century and partly because international communications are now so much more rapid and important. The vocabulary of modern science and technology is very largely international.

THE COMPARATIVE METHOD

The comparative method in historical linguistics is concerned with the reconstruction of an earlier language or

Recon-
struction
of an
earlier
language

earlier state of a language on the basis of a comparison of related words and expressions in different languages or dialects derived from it. The comparative method was developed in the course of the 19th century for the reconstruction of Proto-Indo-European and was subsequently applied to the study of other language families. It depends upon the principle of regular sound change—a principle that, as explained above, met with violent opposition when it was introduced into linguistics by the Neogrammarians in the 1870s but by the end of the century had become part of what might be fairly described as the orthodox approach to historical linguistics. Changes in the phonological systems of languages through time were accounted for in terms of sound laws.

Grimm's law. The most famous of the sound laws is Grimm's law (though Grimm himself did not use the term law). Some of the correspondences accounted for by Grimm's law are given in Table 1. It will be observed that,

Greek	Latin	Gothic	Sanskrit	Slavic
p	p	f	p	p
b	b	p	b	b
ph	f/b	b	bh	b
t	t	θ	t	t
d	d	t	d	d
th	f/d	d	dh	d

when other Indo-European languages, including Latin and Greek, have a voiced unaspirated stop (*b, d*), Gothic has the corresponding voiceless unaspirated stop (*p, t*) and that, when other Indo-European languages have a voiceless unaspirated stop, Gothic has a voiceless fricative (*f, θ*). The simplest explanation would seem to be that, under the operation of what is now called Grimm's law, in some prehistoric period of Germanic (before the development of a number of distinct Germanic languages), the voiced stops inherited from Proto-Indo-European became voiceless and the voiceless stops became fricatives. The situation with respect to the sounds corresponding to the Germanic voiced stops is more complex. Here there is considerable disagreement between the other languages: Greek has voiceless aspirates (*ph, th*), Sanskrit has voiced aspirates (*bh, dh*), Latin has voiceless fricatives in word-initial position (*f*) and voiced stops in medial position (*b, d*), Slavic has voiced stops (*b, d*), and so on. The generally accepted hypothesis is that the Proto-Indo-European sounds from which the Germanic voiced stops developed were voiced aspirates and that they are preserved in Sanskrit but were changed in the other Indo-European languages by the loss of either voice or aspiration. (Latin, having lost the voice in initial position, subsequently changed both of the resultant voiceless aspirates into the fricative *f*, and it lost the aspiration in medial position.) It is easy to see that this hypothesis yields a simpler account of the correspondences than any of the alternatives. It is also in accord with the fact that voiced aspirates are rare in the languages of the world and, unless they are supported by the coexistence in the same language of phonologically distinct voiceless aspirates (as they are in Hindi and other north Indian languages), appear to be inherently unstable.

Proto-Indo-European reconstruction. Reconstruction of the Proto-Indo-European labial stops (made with the lips) and dental stops (made with the tip of the tongue touching the teeth) is fairly straightforward. More controversial is the reconstruction of the Proto-Indo-European sounds underlying the correspondences shown in Table 2. Accord-

Greek	Latin	Gothic	Sanskrit	Slavic
k	k	h	ś	s
g	g	k	j	z
kh	h/g/f	g	h	z
p/t/k	qu	wh	k	k
b/d/g	v/gu	q	g	g
ph/th/kh	f/v/gu	w	gh	g

ing to the most generally accepted hypothesis, there were in Proto-Indo-European at least two distinct series of velar (or "guttural") consonants: simple velars (or palatals), symbolized as **k, *g*, and **gh*, and labiovelars, symbolized as **k^w, *g^w*, and **g^wh*. The labiovelars may be thought of as velar stops articulated with simultaneous lip-rounding. In one group of languages, the labial component is assumed to have been lost, in another group the velar component; and it is only in the Latin reflex of the voiceless **k^w* that both labiality and velarity are retained (*cf.* Latin *quis* from **k^wi-*). It is notable that the languages that have a velar for the Proto-Indo-European labiovelar stops (*e.g.*, Sanskrit and Slavic) have a sibilant or palatal sound (*s* or *ś*) for the Proto-Indo-European simple velars. Earlier scholars attached great significance to this fact and thought that it represented a fundamental division of the Indo-European family into a western and an eastern group. The western group—comprising Celtic, Germanic, Italic, and Greek—is commonly referred to as the centum group; the eastern group—comprising Sanskrit, Iranian, Slavic, and others—is called the satem (*satəm*) group. (The words centum and satem come from Latin and Iranian, respectively, and mean "hundred." They exemplify, with their initial consonant, the two different treatments of the Proto-Indo-European simple velars.) Nowadays less importance is attached to the centum-satem distinction. But it is still generally held that in an early period of Indo-European, there was a sound law operative in the dialect or dialects from which Sanskrit, Iranian, Slavic and the other so-called satem languages developed that had the effect of palatalizing the original Proto-Indo-European velars and eventually converting them to sibilants.

Steps in the comparative method. The information given in the previous paragraphs is intended to illustrate what is meant by a sound law and to indicate the kind of considerations that are taken into account in the application of the comparative method. The first step is to find sets of cognate or putatively cognate forms in the languages or dialects being compared: for example, Latin *decem* = Greek *deka* = Sanskrit *daśa* = Gothic *taihum*, all meaning "ten." From sets of cognate forms such as these, sets of phonological correspondences can be extracted; *e.g.*, (1) Latin *d* = Greek *d* = Sanskrit *d* = Gothic *t*; (2) Latin *e* = Greek *e* = Sanskrit *a* = Gothic *ai* (in the Gothic orthography this represents an *e* sound); (3) Latin *c* (*i.e.*, a *k* sound) = Greek *k* = Sanskrit *ś* = Gothic *h*; (4) Latin *em* = Greek *a* = Sanskrit *a* = Gothic *un*. A set of "reconstructed" phonemes can be postulated (marked with an asterisk by the standard convention) to which the phonemes in the attested languages can be systematically related by means of sound laws. The reconstructed Proto-Indo-European word for "ten" is **dekm*. From this form the Latin word can be derived by means of a single sound change, **m* changes to *em* (usually symbolized as **m > em*); the Greek by means of the sound change **m > a* (*i.e.*, vocalization of the syllabic nasal and loss of nasality); the Sanskrit by means of the palatalizing sound law, **k > ś* and the sound change **m > a* (whether this is assumed to be independent of the law operative in Greek or not); and the Gothic by means of Grimm's law (**d > t, *k > h*) and the sound change **m > un*.

Most 19th-century linguists took it for granted that they were reconstructing the actual word forms of some earlier language, that **dekm*, for example, was a pronounceable Proto-Indo-European word. Many of their successors have been more skeptical about the phonetic reality of reconstructed starred forms like **dekm*. They have said that they are no more than formulae summarizing the correspondences observed to hold between attested forms in particular languages and that they are, in principle, unpronounceable. From this point of view, it would be a matter of arbitrary decision which letter is used to refer to the correspondences: Latin *d* = Greek *d* = Sanskrit *d* = Gothic *t*, and so on. Any symbol would do, provided that a distinct symbol is used for each distinct set of correspondences. The difficulty with this view of reconstruction is that it seems to deny the very *raison d'être* of historical and comparative linguistics. Linguists want to know, if possible, not only that Latin *decem*, Greek *deka*, and so on are related,

Contro-
versy
concerning
Proto-
Indo-
European
velar
consonants

Phono-
logical
correspon-
dences in
cognate
forms

but also the nature of their historical relationship—how they have developed from common ancestral form. They also wish to construct, if feasible, some general theory of sound change. This can be done only if some kind of phonetic interpretation can be given to the starred forms. The important point is that the confidence with which a phonetic interpretation is assigned to the phonemes that are reconstructed will vary from one phoneme to another. It should be clear from the discussion above, for example, that the interpretation of **d* as a voiced dental or alveolar stop is more certain than the interpretation of **k* as a voiceless velar stop. The starred forms are not all on an equal footing from a phonetic point of view.

Criticisms of the comparative method. One of the criticisms directed against the comparative method is that it is based upon a misleading genealogical metaphor. In the mid-19th century, the German linguist August Schleicher introduced into comparative linguistics the model of the "family tree." There is obviously no point in time at which it can be said that new languages are "born" of a common parent language. Nor is it normally the case that the parent language "lives on" for a while, relatively unchanged, and then "dies." It is easy enough to recognize the inappropriateness of these biological expressions. No less misleading, however, is the assumption that languages descended from the same parent language will necessarily diverge, never to converge again, through time. This assumption is built into the comparative method as it is traditionally applied. And yet there are many clear cases of convergence in the development of well-documented languages. The dialects of England are fast disappearing and are far more similar in grammar and vocabulary today than they were even a generation ago. They have been strongly influenced by the standard language. The same phenomenon, the replacement of nonstandard or less prestigious forms with forms borrowed from the standard language or dialect, has taken place in many different places at many different times. It would seem, therefore, that one must reckon with both divergence and convergence in the diachronic development of languages: divergence when contact between two speech communities is reduced or broken and convergence when the two speech communities remain in contact and when one is politically or culturally dominant.

The comparative method presupposes linguistically uniform speech communities and independent development after sudden, sharp cleavage. Critics of the comparative method have pointed out that this situation does not generally hold. In 1872 a German scholar, Johannes Schmidt, criticized the family-tree theory and proposed instead what is referred to as the wave theory, according to which different linguistic changes will spread, like waves, from a politically, commercially, or culturally important centre along the main lines of communication, but successive innovations will not necessarily cover exactly the same area. Consequently, there will be no sharp distinction between contiguous dialects, but, in general, the further apart two speech communities are, the more linguistic features there will be that distinguish them (see below *Dialectology and linguistic geography*).

Internal reconstruction. The comparative method is used to reconstruct earlier forms of a language by drawing upon the evidence provided by other related languages. It may be supplemented by what is called the method of internal reconstruction. This is based upon the existence of anomalous or irregular patterns of formation and the assumption that they must have developed, usually by sound change, from earlier regular patterns. For example, the existence of such patterns in early Latin as *honoris* ("honor" : "of honor") and others in contrast with *orator* : *oratoris* ("orator" : "of the orator") and others might lead to the supposition that *honoris* developed from an earlier **honosis*. In this case, the evidence of other languages shows that **s* became *r* between vowels in an earlier period of Latin. But it would have been possible to reconstruct the earlier intervocalic **v* with a fair degree of confidence on the basis of the internal evidence alone. Clearly, internal reconstruction depends upon the structural approach to linguistics.

The most recent development in the field of histori-

cal and comparative linguistics has come from the theory of generative grammar (see above *Transformational-generative grammar*). If the grammar and phonology of a language are described from a synchronic point of view as an integrated system of rules, then the grammatical and phonological similarities and differences between two closely related languages, or dialects, or between two diachronically distinct states of the same language can be described in terms of the similarities and differences in two descriptive rule systems. One system may contain a rule that the other lacks (or may restrict its application more or less narrowly); one system may differ from the other in that the same set of rules will apply in a different order in the one system from the order in which they apply in the other. Language change may thus be accounted for in terms of changes introduced into the underlying system of phonological and grammatical rules (including the addition, loss, or reordering of rules) during the process of language acquisition. So far these principles have been applied principally to sound change. There has also been a little work done on diachronic syntax.

LANGUAGE CLASSIFICATION

There are two kinds of classification of languages practiced in linguistics: genetic (or genealogical) and typological. The purpose of genetic classification is to group languages into families according to their degree of diachronic relatedness. For example, within the Indo-European family, such subfamilies as Germanic or Celtic are recognized; these subfamilies comprise German, English, Dutch, Swedish, Norwegian, Danish, and others, on the one hand, and Irish, Welsh, Breton, and others, on the other. So far, most of the languages of the world have been grouped only tentatively into families, and many of the classificatory schemes that have been proposed will no doubt be radically revised as further progress is made.

A typological classification groups languages into types according to their structural characteristics. The most famous typological classification is probably that of isolating, agglutinating, and inflecting (or fusional) languages, which was frequently invoked in the 19th century in support of an evolutionary theory of language development. Roughly speaking, an isolating language is one in which all the words are morphologically unanalyzable (*i.e.*, in which each word is composed of a single morph); Chinese and, even more strikingly, Vietnamese are highly isolating. An agglutinating language (*e.g.*, Turkish) is one in which the word forms can be segmented into morphs, each of which represents a single grammatical category. An inflecting language is one in which there is no one-to-one correspondence between particular word segments and particular grammatical categories. The older Indo-European languages tend to be inflecting in this sense. For example, the Latin suffix *-is* represents the combination of categories "singular" and "genitive" in the word form *hominis* "of the man," but one part of the suffix cannot be assigned to "singular" and another to "genitive," and *-is* is only one of many suffixes that in different classes (or declensions) of words represent the combination of "singular" and "genitive."

There is, in principle, no limit to the variety of ways in which languages can be grouped typologically. One can distinguish languages with a relatively rich phonemic inventory from languages with a relatively poor phonemic inventory, languages with a high ratio of consonants to vowels from languages with a low ratio of consonants to vowels, languages with a fixed word order from languages with a free word order, prefixing languages from suffixing languages, and so on. The problem lies in deciding what significance should be attached to particular typological characteristics. Although there is, not surprisingly, a tendency for genetically related languages to be typologically similar in many ways, typological similarity of itself is no proof of genetic relationship. Nor does it appear true that languages of a particular type will be associated with cultures of a particular type or at a certain stage of development. What has emerged from recent work in typology is that certain logically unconnected features tend to occur together, so that the presence of feature A in a given

Model of
the family
tree

Use of
irregular
forms for
recon-
struction

Typologi-
cal classifi-
cations

language will tend to imply the presence of feature B. The discovery of unexpected implications of this kind calls for an explanation and gives a stimulus to research in many branches of linguistics.

Linguistics and other disciplines

PSYCHOLINGUISTICS

The term psycholinguistics was coined in the 1940s and came into more general use after the publication of Charles E. Osgood and Thomas A. Sebeok's *Psycholinguistics: A Survey of Theory and Research Problems* (1954), which reported the proceedings of a seminar sponsored in the United States by the Social Science Research Council's Committee on Linguistics and Psychology.

The boundary between linguistics (in the narrower sense of the term: see the introduction of this article) and psycholinguistics is difficult, perhaps impossible, to draw. So too is the boundary between psycholinguistics and psychology. What characterizes psycholinguistics as it is practiced today as a more or less distinguishable field of research is its concentration upon a certain set of topics connected with language and its bringing to bear upon them the findings and theoretical principles of both linguistics and psychology. The range of topics that would be generally held to fall within the field of psycholinguistics nowadays is rather narrower, however, than that covered in the survey by Osgood and Sebeok.

Language acquisition by children. One of the topics most central to psycholinguistic research is the acquisition of language by children. The term "acquisition" is preferred to "learning," because "learning" tends to be used by psychologists in a narrowly technical sense, and many psycholinguists believe that no psychological theory of learning, as currently formulated, is capable of accounting for the process whereby children, in a relatively short time, come to achieve a fluent control of their native language. Since the beginning of the 1960s, research on language acquisition has been strongly influenced by Chomsky's theory of generative grammar, and the main problem to which it has addressed itself has been how it is possible for young children to infer the grammatical rules underlying the speech they hear and then to use these rules for the construction of utterances that they have never heard before. It is Chomsky's conviction, shared by a number of psycholinguists, that children are born with a knowledge of the formal principles that determine the grammatical structure of all languages, and that it is this innate knowledge that explains the success and speed of language acquisition. Others have argued that it is not grammatical competence as such that is innate but more general cognitive principles and that the application of these to language utterances in particular situations ultimately yields grammatical competence. Many recent works have stressed that all children go through the same stages of language development regardless of the language they are acquiring. It has also been asserted that the same basic semantic categories and grammatical functions can be found in the earliest speech of children in a number of different languages operating in quite different cultures in various parts of the world.

Although Chomsky was careful to stress in his earliest writings that generative grammar does not provide a model for the production or reception of language utterances, there has been a good deal of psycholinguistic research directed toward validating the psychological reality of the units and processes postulated by generative grammarians in their descriptions of languages. Experimental work in the early 1960s appeared to show that nonkernel sentences took longer to process than kernel sentences and, even more interestingly, that the processing time increased proportionately with the number of optional transformations involved. More recent work has cast doubt on these findings, and most psycholinguists are now more cautious about using grammars produced by linguists as models of language processing. Nevertheless, generative grammar continues to be a valuable source of psycholinguistic experimentation, and the formal properties of language, discovered or more adequately discussed

by generative grammarians than they have been by others, are generally recognized to have important implications for the investigation of short-term and long-term memory and perceptual strategies.

Speech perception. Another important area of psycholinguistic research that has been strongly influenced by recent theoretical advances in linguistics and, more especially, by the development of generative grammar is speech perception. It has long been realized that the identification of speech sounds and of the word forms composed of them depends upon the context in which they occur and upon the hearer's having mastered, usually as a child, the appropriate phonological and grammatical system. Throughout the 1950s, work on speech perception was dominated (as was psycholinguistics in general) by information theory, according to which the occurrence of each sound in a word and each word in an utterance is statistically determined by the preceding sounds and words. Information theory is no longer as generally accepted as it was a few years ago, and more recent research has shown that in speech perception the cues provided by the acoustic input are interpreted, unconsciously and very rapidly, with reference not only to the phonological structure of the language but also to the more abstract levels of grammatical organization.

Other areas of research. Other areas of psycholinguistics that should be briefly mentioned are the study of aphasia and neurolinguistics. The term aphasia is used to refer to various kinds of language disorders; recent work has sought to relate these, on the one hand, to particular kinds of brain injury and, on the other, to psychological theories of the storage and processing of different kinds of linguistic information. One linguist has put forward the theory that the most basic distinctions in language are those that are acquired first by children and are subsequently most resistant to disruption and loss in aphasia. This, though not disproved, is still regarded as controversial. Two kinds of aphasia are commonly distinguished. In motor aphasia the patient manifests difficulty in the articulation of speech or in writing and may produce utterances with a simplified grammatical structure, but his comprehension is not affected. In sensory aphasia the patient's fluency may be unaffected, but his comprehension will be impaired and his utterances will often be incoherent.

Neurolinguistics should perhaps be regarded as an independent field of research rather than as part of psycholinguistics. In 1864 it was shown that motor aphasia is produced by lesions in the third frontal convolution of the left hemisphere of the brain. Shortly after the connection had been established between motor aphasia and damage to this area (known as Broca's area), the source of sensory aphasia was localized in lesions of the posterior part of the left temporal lobe. More recent work has confirmed these findings. The technique of electrically stimulating the cortex in conscious patients has enabled brain surgeons to induce temporary aphasia and so to identify a "speech area" in the brain. It is no longer generally believed that there are highly specialized "centres" within the speech area, each with its own particular function; but the existence of such a speech area in the dominant hemisphere of the brain (which for most people is the left hemisphere) seems to be well established. The posterior part of this area is involved more in the comprehension of speech and the construction of grammatically and semantically coherent utterances, and the anterior part is concerned with the articulation of speech and with writing. Little is yet known about the operation of the neurological mechanisms underlying the storage and processing of language. (See also the articles entitled PERCEPTION; SPEECH.)

SOCIOLINGUISTICS

Delineation of the field. Just as it is difficult to draw the boundary between linguistics and psycholinguistics and between psychology and psycholinguistics, so it is difficult to distinguish sharply between linguistics and sociolinguistics and between sociolinguistics and sociology. There is the further difficulty that, because the boundary between sociology and anthropology is also unclear, sociolinguistics merges with anthropological linguistics (see below).

Information theory and speech perception

Theories of language acquisition

Speech area in the brain

It is frequently suggested that there is a conflict between the sociolinguistic and the psycholinguistic approach to the study of language, and it is certainly the case that two distinct points of view are discernible in the literature at the present time. Chomsky has described linguistics as a branch of cognitive psychology, and neither he nor most of his followers have yet shown much interest in the relationship between language and its social and cultural matrix. On the other hand, many modern schools of linguistics that have been very much concerned with the role of language in society would tend to relate linguistics more closely to sociology and anthropology than to any other discipline. It would seem that the opposition between the psycholinguistic and the sociolinguistic viewpoint must ultimately be transcended. The acquisition of language, a topic of central concern to psycholinguists, is in part dependent upon and in part itself determines the process of socialization; and the ability to use one's native language correctly in the numerous socially prescribed situations of daily life is as characteristic a feature of linguistic competence, in the broad sense of this term, as is the ability to produce grammatical utterances. Some of the most recent work in sociolinguistics and psycholinguistics has sought to widen the notion of linguistic competence in this way. So far, however, sociolinguistics and psycholinguistics tend to be regarded as relatively independent areas of research.

Language
and social-
ization

Social dimensions. Language is probably the most important instrument of socialization that exists in all human societies and cultures. It is largely by means of language that one generation passes on to the next its myths, laws, customs, and beliefs, and it is largely by means of language that the child comes to appreciate the structure of the society into which he is born and his own place in that society.

As a social force, language serves both to strengthen the links that bind the members of the same group and to differentiate the members of one group from those of another. In many countries there are social dialects as well as regional dialects, so that it is possible to tell from a person's speech not only where he comes from but what class he belongs to. In some instances social dialects can transcend regional dialects. This is notable in England, where standard English in the so-called Received Pronunciation (RP) can be heard from members of the upper class and upper middle class in all parts of the country. The example of England is but an extreme manifestation of a tendency that is found in all countries: there is less regional variation in the speech of the higher than in that of the lower socioeconomic classes. In Britain and the United States and in most of the other English-speaking countries, people will almost always use the same dialect, regional or social, however formal or informal the situation and regardless of whether their listeners speak the same dialect or not. (Relatively minor adjustments of vocabulary may, however, be made: an Englishman speaking to an American may employ the word "elevator" rather than "lift" and so on.) In many communities throughout the world, it is common for members to speak two or more different dialects and to use one dialect rather than another in particular social situations. This is commonly referred to as code-switching. Code-switching may operate between two distinct languages (e.g., Spanish and English among Puerto Ricans in New York) as well as between two dialects of the same language. The term diglossia (rather than bilingualism) is frequently used by sociolinguists to refer to this by no means uncommon phenomenon.

In every situation, what one says and how one says it depends upon the nature of that situation, the social role being played at the time, one's status vis-à-vis that of the person addressed, one's attitude towards him, and so on. Language interacts with nonverbal behaviour in social situations and serves to clarify and reinforce the various roles and relationships important in a particular culture. Sociolinguistics is far from having satisfactorily analyzed or even identified all the factors involved in the selection of one language feature rather than another in particular situations. Among those that have been discussed in relation to various languages are: the formality or informality of the situation; power and solidarity relationships be-

Variables
influencing
language
usage

tween the participants; differences of sex, age, occupation, socioeconomic class, and educational background; and personal or transactional situations. Terms such as style and register (as well as a variety of others) are employed by many linguists to refer to the socially relevant dimensions of phonological, grammatical, and lexical variation within one language. So far there is very little agreement as to the precise application of such terms. (For further treatment of sociolinguistics, see the section *Dialects* in the article LANGUAGE.)

OTHER RELATIONSHIPS

Anthropological linguistics. The fundamental concern of anthropological linguistics is to investigate the relationship between language and culture. To what extent the structure of a particular language is determined by or determines the form and content of the culture with which it is associated remains a controversial question. Vocabulary differences between languages correlate obviously enough with cultural differences, but even here the interdependence of language and culture is not so strong that one can argue from the presence or absence of a corresponding cultural difference. For example, from the fact that English—unlike French, German, Russian, and many other languages—distinguishes lexically between monkeys and apes, one cannot conclude that there is an associated difference in the cultural significance attached to these animals by English-speaking societies. Some of the major grammatical distinctions in certain languages may have originated in culturally important categories (e.g., the distinction between an animate and an inanimate gender). But they seem to endure independently of any continuing cultural significance. The "Whorfian hypothesis" (the thesis that one's thought and even perception are determined by the language one happens to speak), in its strong form at least, is no longer debated as vigorously as it was a few years ago. Anthropologists continue to draw upon linguistics for the assistance it can give them in the analysis of such topics as the structure of kinship. A more recent development, but one that has not so far produced any very substantial results, is the application of notions derived from generative grammar to the analysis of ritual and other kinds of culturally prescribed behaviour.

The
"Whorfian
hypothesis"

Computational linguistics. By computational linguistics is meant no more than the use of electronic digital computers in linguistic research. At a theoretically trivial level, computers are employed to scan texts and to produce, more rapidly and more reliably than was possible in the past, such valuable aids to linguistic and stylistic research as word lists, frequency counts, and concordances. Theoretically more interesting, though much more difficult, is the automatic grammatical analysis of texts by computer. Considerable progress was made in this area by research groups working on machine translation and information retrieval in the United States, Great Britain, the Soviet Union, France, and a few other countries in the decade between the mid-1950s and the mid-1960s. But much of the original impetus for this work disappeared, for a time at least, in part because of the realization that the theoretical problems involved in machine translation are much more difficult than they were at first thought to be and in part as a consequence of a loss of interest among linguists in the development of discovery procedures. Whether automatic syntactic analysis and fully automatic high-quality machine translation are even feasible in principle remains a controversial question.

Use of
comput-
ers in
linguistic
research

Mathematical linguistics. What is commonly referred to as mathematical linguistics comprises two areas of research: the study of the statistical structure of texts and the construction of mathematical models of the phonological and grammatical structure of languages. These two branches of mathematical linguistics, which may be termed statistical and algebraic linguistics, respectively, are typically distinct. Attempts have been made to derive the grammatical rules of languages from the statistical structure of texts written in those languages, but such attempts are generally thought to have been not only unsuccessful so far in practice but also, in principle, doomed to failure. That languages have a statistical structure is a fact well

known to cryptographers. Within linguistics, it is of considerable typological interest to compare languages from a statistical point of view (the ratio of consonants to vowels, of nouns to verbs, and so on). Statistical considerations are also of value in stylistics (see below).

Algebraic
linguistics

Algebraic linguistics derives principally from the work of Noam Chomsky in the field of generative grammar (see above *Chomsky's grammar*). In his earliest work Chomsky described three different models of grammar—finite-state grammar, phrase-structure grammar, and transformational grammar—and compared them in terms of their capacity to generate all and only the sentences of natural languages and, in doing so, to reflect in an intuitively satisfying manner the underlying formal principles and processes. Other models have also been investigated, and it has been shown that certain different models are equivalent in generative power to phrase-structure grammars. The problem is to construct a model that has all the formal properties required to handle the processes found to be operative in languages but that prohibits rules that are not required for linguistic description. It is an open question whether such a model, or one that approximates more closely to this ideal than current models do, will be a transformational grammar or a grammar of some radically different character.

Stylistics. The term stylistics is employed in a variety of senses by different linguists. In its widest interpretation it is understood to deal with every kind of synchronic variation in language other than what can be ascribed to differences of regional dialect. At its narrowest interpretation it refers to the linguistic analysis of literary texts. One of the aims of stylistics in this sense is to identify those features of a text that give it its individual stamp and mark it as the work of a particular author. Another is to identify the linguistic features of the text that produce a certain aesthetic response in the reader. The aims of stylistics are the traditional aims of literary criticism. What distinguishes stylistics as a branch of linguistics (for those who regard it as such) is the fact that it draws upon the methodological and theoretical principles of modern linguistics.

Philosophy of language. The analysis of language has always been a subject of particular concern to philosophers, and traditional grammar was strongly influenced by the dominant philosophical attitudes of the day. Modern linguistics and modern philosophical theories have so far had little influence on one another. Some philosophers have shown an interest in Chomsky's controversial suggestion that work in generative grammar lends support to the rationalists in their long-standing dispute about the source of human knowledge. Potentially more fruitful, perhaps, is the interest shown by a number of linguists in philosophical treatments of reference, quantification, and presupposition, in systems of modal logic, and in the work of the so-called philosophers of ordinary language.

Applied linguistics. In the sense in which the term applied linguistics is most commonly used nowadays it is restricted to the application of linguistics to language teaching. Much of the recent expansion of linguistics as a subject of teaching and research in the universities in many countries has come about because of its value, actual and potential, for writing better language textbooks and devising more efficient methods of teaching languages. Linguistics is also widely held to be relevant to the training of teachers of the deaf and speech therapists. Outside the field of education in the narrower sense, applied linguistics (and, more particularly, applied sociolinguistics) has an important part to play in what is called language planning; *i.e.*, in advising governments, especially in recent created states, as to which language or dialect should be made the official language of the country and how it should be standardized. (J.Lyo./Ed.)

Dialectology and linguistic geography

DIALECT GEOGRAPHY

Dialect study as a discipline—dialectology—dates from the first half of the 19th century, when local dialect dictionaries and dialect grammars first appeared in western

Europe. Soon thereafter, dialect maps were developed; most often they depicted the division of a language's territory into regional dialects. The 19th-century rise of nationalism, coupled with the Romantic view of dialects and folklore as manifestations of the ethnic soul, furnished a great impetus for dialectology.

19th-
century
impetus
for dialect
study

Early dialect studies. The first dialect dictionaries and grammars were most often written by scholars describing the dialect of their birthplace or by fieldworkers whose main method of investigation was free conversation with speakers of the dialect, usually older persons and, preferably, those who showed the least degree of literacy and who had travelled as little as possible. Many of these grammars and dictionaries recorded dialectal traits that deviated from the standard language. In the second half of the 19th century, when historical and comparative linguistic study was flourishing, it became customary to focus attention on the fate of particular elements of the archaic language in a given dialect; *e.g.*, the changes that Latin vowels and consonants underwent when used in different positions in a particular Romance dialect.

With the accumulation of dialectal data, investigators became increasingly conscious of the inadequacy of viewing dialects as internally consistent units that were sharply differentiated from neighbouring dialects. It became more and more clear that each dialectal element or phenomenon refused to stay neatly within the borders of a single dialect area and that each had its own isogloss; consequently, maps of dialects would have to be replaced by maps showing the distribution of each particular feature. While sound scientifically, the preparation and compilation of such maps, called linguistic atlases, is a difficult, costly, voluminous, and time-consuming job.

Dialect atlases. Dialect atlases are compiled on the basis of investigations of the dialects of a large number of places; a questionnaire provides uniform data. There are two basic methods of data collection: fieldwork and survey by correspondence. Fieldwork, in which a trained investigator transcribes dialectal forms directly (or on tape), affords more precise data and enables the questionnaire to include a greater number of diverse questions; but it implies a necessarily limited number of points to be covered. The advantage of the correspondence method lies in its ability to encompass more points at less cost and with less time expended in gathering the data. On the other hand, rural schoolteachers, normally the persons who complete such questionnaires, can answer only a relatively small number of questions and often imperfectly.

The first large-scale enterprise in linguistic geography was the preparation of the German linguistic atlas. In the 1880s, the initiator of this great undertaking, Georg Wenker, composed 40 test sentences that illustrated most of the important ways in which dialects differed and sent them to schoolmasters in over 40,000 places in the German Empire. The sentences were to be translated into the local dialect. Publication of the results was not begun until 1926; the main cause of the delay was the enormous quantity of material to be arranged and analyzed.

The famous French linguistic atlas of Jules Gilliéron and Edmond Edmont was based on a completely different concept. Using a questionnaire of about 2,000 words and phrases that Gilliéron had composed, Edmont surveyed 639 points in the French-speaking area. The atlas, compiled under the direction of Gilliéron, was published in fascicles from 1902 to 1912 and furnished both a strong stimulus and the basic model for work on linguistic atlases elsewhere in the world. European linguists, especially in Romance- and Germanic-speaking countries, were the first to participate in such atlas projects. One of the most significant contributions is the linguistic atlas of Italy and southern Switzerland by Karl Jaberg and Jakob Jud; it appeared from 1928 to 1940. Particularly noteworthy in its attention to precise definitions of meaning, this atlas often used illustrations and described objects and actions of village life denoted by the questionnaire's words.

The
French
atlas

At present, dialects of virtually all European languages have been treated in linguistic geography studies. In some countries, data are still being collected and classified and maps are being drawn, but in others a second generation

Language
teaching

of atlases is already under way. French dialectologists, for instance, are now working on regional atlases that will complement data contained in the *Atlas linguistique de la France*. In England, work began in 1946, under the direction of Harold Orton and Eugene Dieth: the first volume of the *Survey of English Dialects* was published in 1962. In Slavic-speaking countries, work is now under way both on atlases of separate Slavic languages and on the large general Slavic linguistic atlas that will cover nearly 1,000 locales in all parts of European territory where Slavic languages are spoken. Outside Europe, the greatest amount of work in linguistic geography has been completed in Japan and in the United States.

As early as 1905–06, a committee of Japanese dialectologists published the first linguistic atlas of Japan in two volumes, one devoted to phonology and one to morphology. Subsequent work has been done on a new atlas of Japan as a whole and on several regional atlases. The extensive activity of Chinese specialists has concentrated on descriptions of particular local and regional dialects. The Chinese situation is a peculiar one because of the enormous number of people who speak Chinese, the very significant dialectal differentiation (certain dialects, particularly those in the south of China, would be considered by Western standards as separate languages), and the nature of the Chinese script. Chinese characters do not represent sounds but concepts. Because of this, the written language can be read without difficulty in many different dialect areas, although its spoken form varies greatly from one region to another.

Because of the enormous size of the United States, atlas surveys were done by region. Between 1931 and 1933, fieldworkers under the direction of the linguist Hans Kurath surveyed 213 New England communities; the results were published in the *Linguistic Atlas of New England* (with 734 maps) in 1939–43. Based on the methodological experience of Jaberg and Jud in their atlas of Italy and southern Switzerland, this work involved systematic investigations not only among the relatively uneducated but also among better educated, more cultured informants and among the very well educated, cultured, and informed members of a community. Thus the dimension of social stratification of language was introduced into linguistic geography, and valuable material about regional linguistic standards became available.

After 1933, fieldwork was extended to the other Atlantic states. Lack of financial support, however, has hindered the publication of these atlases. Nevertheless, several works based on the material gathered have appeared, among them Kurath's *Word Geography of the Eastern United States*, E. Bagby Atwood's *Survey of Verb Forms in the Eastern United States*, and Kurath's and McDavid's *Pronunciation of English in the Atlantic States*. Independent work was carried out in other U.S. regions, mainly with an adapted form of the questionnaire developed for the Atlantic states; only introductions or summaries of material in the files have been published, however, because of lack of funds.

The most effective and thorough—as well as the most expensive—way of presenting data in linguistic atlases is by printing the actual responses to questionnaire items right on the maps. Phenomena of linguistic geography, however, are usually represented by geometric symbols or figures at the proper points on the map or, even more summarily, by the drawing of isoglosses (linguistic boundaries) or by shading or colouring the areas of particular features.

Only dialect atlases can furnish the complexity of data of the major dialectal phenomena in a multitude of geographic locations in a manner that both assures commensurability of the data and allows a panoramic examination of the whole gamut of data. The inventory of linguistic phenomena is so rich, however, that no one questionnaire can encompass it all. Moreover, the use of a questionnaire unavoidably brings about a schematization of answers that is lacking in spontaneity. For these reasons, other kinds of publications, such as dialect dictionaries or monographs based on extensive free conversation with speakers of local dialects, are indispensable complements to linguistic atlases.

The value and applications of dialectology. The scientific interest of dialectology lies in the fact that dialects are a valuable source of information about popular culture. They reflect not only the history of a language but, to a great extent, the ethnic, cultural, and even political history of a people as well. A knowledge of dialectal facts provides practical guidance to school systems that are trying to teach the standard language to an ever greater number of pupils.

In the 1930s the value of dialectology to the study of language types became apparent. Because dialects greatly outnumber standard languages, they provide a much greater variety of phenomena than languages and thus have become the main source of information about the types of phenomena possible in linguistic systems. Also, in some languages, but not in others, an extremely wide structural variation among dialects has been found. In the Balkan region, where two closely related Slavic languages, Serbo-Croatian and Slovene, are spoken, dialects are found with synthetic declension (case endings, as in Latin) and analytic declension (use of prepositions and word order, as in English). In addition, there are among these dialects complex systems of verbal tenses contrasting with simple ones, as well as dialects with or without the dual number or the neuter gender. The dialects of Serbo-Croatian and Slovene also exhibit almost every type of prosodic structure (*e.g.*, tone, stress, length) found in European languages. Some dialects differentiate long and short vowels or rising and falling accents, while others do not; and in some, but not all, of them stress fulfills a grammatical function. Of the several dozen vowel and diphthong sounds that occur in these dialects, only five are common to all of them; all the rest are restricted to relatively small areas. All of this rich variety contrasts sharply with the relative structural uniformity of the English language—not only in the United States but wherever it is spoken. (The outstanding exceptions are the creolized dialects, which are distinguished by far-reaching structural peculiarities.)

SOCIAL DIALECTOLOGY

The methodology of generative grammar was first applied to dialectology in the 1960s, when the use of statistical means to measure the similarity or difference between dialects also became increasingly common. The most important development of that time, however, was the rapid growth of methods for investigating the social variation of dialects; social variation, in contrast to geographic variation, is prominent in the United States, above all in large urban centres. In cities such as New York, a whole scale of speech variation can be found to correlate with the social status and educational level of the speakers. In addition, age groups exhibit different patterns, but such patterns of variation differ from one social stratum to another. Still another dimension of variation, especially important in the United States, is connected with the race and ethnic origin of a speaker as well as with the speaker's date of immigration. So-called Black English has been influenced by the southeastern U.S. origin of most of the black population of non-southern U.S. regions: many Black English peculiarities are in reality transplanted southeastern dialectal traits.

Normally, speakers of one of the social dialects of a city possess at least some awareness of the other dialects. In this way, speech characteristics also become subjectively integrated into the system of signs indicating social status. And, in seeking to enhance their social status, poorer and less educated speakers may try to acquire the dialect of the socially prestigious. Certain groups—*e.g.*, blacks and the working class—however, will, under certain conditions, show a consciousness of solidarity and a tendency to reject members who imitate either the speech or other types of behaviour of models outside their own social group.

As a consequence of an individual's daily contacts with speakers of the various social dialects of a city, elements of the other dialects are imperceptibly drawn into his dialect. The collective result of such experiences is the spread of linguistic variables—*i.e.*, groups of variants (sounds or grammatical phenomena) primarily determined by social (educational, racial, age, class) influences, an example be-

Variations in Serbo-Croatian and Slovene dialects

Social correlations with speech variation

Works on American English

ing the existence of the two forms "He don't know" and the standard "He doesn't know." Traits representing variables in intergroup relations can become variable features in the speech of individuals as well; *i.e.*, an individual may employ two or more variants for the same feature in his own speech, such as "seeing" and "seein'" or "he don't" and "he doesn't." The frequency of usage for each variable varies with the individual speaker as well as with the social group. There are intermediate stages of frequency between different social groups and entire scales of transitions between different age groups, thus creating even greater variation within the dialect of an individual. The variables also behave differently in the various styles of written or spoken language used by each speaker.

The study of variables is one of the central tasks of any investigation of the dialects of American cities. Applying the statistical methods of modern sociology, linguists have worked out investigative procedures sharply different from those of traditional dialectology. The chief contributor has been William Labov, the pioneer of social dialectology in the U.S. The basic task is to determine the correlation between a group of linguistic variables—such as the different ways of pronouncing a certain vowel—and extralinguistic variables, such as education, social status, age, and race. For a reasonable degree of statistical reliability, one must record a great number of speakers. In general, several examples of the same variable must be elicited from each individual in order to examine the frequency and probability of its usage. Accordingly, the number of linguistic variables that can be examined is quite limited, in comparison with the number of dialectal features normally recorded by traditional fieldworkers in rural communities; in these situations, the investigator is often satisfied with one or two responses for each feature.

A completely new, flexible, and imaginative method of interviewing is needed for such work in urban centres, as well as new ways of finding and making contact with informants. One example is Labov's method for testing the fate of final and preconsonantal *r* in speakers of different social levels. Choosing three New York City department stores, each oriented to a completely different social stratum, he approached a large number of salesladies, asking each of them about the location of a certain department that he knew to be on the fourth floor. Thus, their answers always contained two words with potential *r*'s—"fourth" and "floor." This shortcut enabled Labov to establish in a relatively short time that the salesladies in the store with richer customers clearly tended to use "*r*-full" forms, whereas those in the stores geared to the poorer social strata more commonly used "*r*-less" forms.

Social dialectology has focused on the subjective evaluation of linguistic features and the degree of an individual's linguistic security, phenomena that have considerable influence on linguistic change. Linguistic scientists, in studying the mechanism of such change, have found that it seems to proceed gradually from one social group to another, always attaining greater frequency among the young. Social dialectology also has great relevance for a society as a whole, in that the data it furnishes will help deal with the extremely complex problems connected with the speech of the socially underprivileged, especially of minority groups. Thus, the recent emphasis on the speech of minority groups, such as the Black English of American cities, is not a chance phenomenon. Specific methods for such investigation are being developed, as well as ways of applying the results of such investigation to educational policies.

(P.I./Ed.)

Semantics

Semantics is the philosophical and scientific study of meaning. The term is one of a group of English words formed from the various derivatives of the Greek verb *sēmainō* ("to mean" or "to signify"). The noun semantics and the adjective semantic are derived from *sēmantikos* ("significant"); semiotic (adjective and noun) comes from *sēmeiōtikos* ("pertaining to signs"); semology from *sēma* ("sign") + *logos* ("account"); and semasiology from *sēmasia* ("signification") + *logos* ("account"). It is difficult

to formulate a distinct definition for each of these terms because their use largely overlaps in the literature despite individual preferences. Semantics is a relatively new field of study, and its originators, often working independently of one another, felt the need to coin a new name for the new discipline; hence the variety of terms denoting the same subject. The word semantics has ultimately prevailed as a name for the doctrine of meaning, in particular, of linguistic meaning. Semiotic is still used, however, to denote a broader field: the study of sign-using behaviour in general.

MODERN DEVELOPMENT OF SEMANTICS

The concern with meaning, always present for philosophers and linguists, greatly increased in the decades following World War II. The sudden rise of interest in meaning can be attributed to an interaction of several lines of development in various disciplines. From the middle of the 19th century onward, logic, the formal study of reasoning, underwent a period of growth unparalleled since the time of Aristotle. Although the main motivation for a renewed interest in logic was a search for the foundations of mathematics, the chief protagonists of this effort—notably the German mathematician Gottlob Frege and the English philosopher Bertrand Russell—extended their inquiry into the domain of the natural languages, which are the original media of human reasoning. The influence of mathematical thinking, and of mathematical logic in particular, however, left a permanent mark on the subsequent study of semantics.

Positivist theory. This mark is nowhere more obvious than in the semantic theories offered by the Neopositivists of the Vienna Circle, which flourished in the 1920s and 1930s, and which was composed of philosophers, mathematicians, and scientists who discussed the methodology and epistemology of science. To such "logical" Positivists as the German-born philosopher Rudolf Carnap, for instance, the symbolism of modern logic represented the grammar (syntax) of an "ideal" language. Because the Logical Positivists were, at the same time, radical Empiricists (observationalists) in their philosophy, the semantics of their ideal language has been given in terms of a tie connecting the symbols of this language with observable entities in the world, or the data of one's sense experience, or both. Against such a rigid ideal as logic, natural language appeared to these philosophers as something primitive, vague, inaccurate, and confused. Moreover, since a large part of ordinary and philosophical discourse, particularly that concerning metaphysical and moral issues, could not be captured by the ideal language, the Positivist approach provided a way to brand all such talk as nonsensical, or at least as "cognitively" meaningless. Accordingly, the Positivists engaged in a prolonged, and largely unsuccessful, effort to formulate a criterion of meaningfulness in terms of empirical verifiability with respect to the sentences formed in natural language.

Whorfian views. Another source of dissatisfaction with the vernacular was made apparent shortly before World War II by the work of the American anthropological linguist Benjamin Lee Whorf. Whorf's famous thesis of linguistic relativity implied that the particular language a person learns and uses determines the framework of his perception and thought. If that language is vague and inaccurate, as the Positivists suggested, or is burdened with the prejudices and superstitions of an ignorant past, as some cultural anthropologists averred, then it is bound to render the user's thinking—and his mental life itself—confused, prejudiced, and superstitious. The Polish-American semanticist Alfred Korzybski, the founder of the movement called General Semantics, believed that the cure for such vague and superstition-laden language lay in a radical revision of linguistic habits in the light of modern science.

School of natural language. Natural language did not remain without champions in the face of this combined onslaught from the logicians and the Whorfians. A reaction started in England, first in Cambridge, then in Oxford. Influenced by the English philosopher George Edward Moore but more so by the "converted" Vienna-born Positivist Ludwig Wittgenstein, the philosophy of "ordinary

Semantic theories of the Logical Positivists

Application of modern statistical methods to dialectology

language" (also known as the Oxford school) came into its own in the 1940s. According to the philosophers of this group, natural language, far from being the crude instrument the Positivists alleged it to be, provides the basic and unavoidable matrix of all thought, including philosophical reflections. Any "ideal" language, therefore, can make sense only as a parasitical extension of, and never as a substitute for, the natural language. Philosophical problems arise as a result of a failure to see the workings of man's language; they are bound to "dissolve" with improved understanding. These assumptions provided a mighty impetus to reflect upon the vernacular language, including its minute points of grammar and fine nuances of meaning. Indeed, some of the later representatives of this approach, particularly the English philosopher John L. Austin, became renowned as much among linguists as among philosophers.

Modern grammatical influences. In the 1950s the science of linguistics itself rose to the challenges that had been coming chiefly from philosophical quarters. The development of transformational, or generative, grammar, initiated by the work of the U.S. linguists Zellig S. Harris and Noam Chomsky, opened a deeper insight into the syntax of the natural languages. Instead of merely providing a structural description (parsing) of sentences, this approach demonstrates how sentences are built up, step by step, from some basic ingredients. In the hands of the philosopher, this powerful new grammar not only served to counter the positivistic charge of imprecision laid against natural language but aided him in his own work of conceptual clarification. Moreover, the generative approach promised further results: since the late 1960s some steps have been taken to develop a generative semantics for natural languages, in addition to a generative syntax.

PHILOSOPHICAL VIEWS ON MEANING

Meaning and reference. On a rather unsophisticated level the problem of meaning can be approached through the following steps. The perception of certain physical entities (objects, marks, sounds, and so on) might lead an intelligent being to the thought of another thing with some regularity. For example, the sight of smoke evokes the idea of fire, footprints on the sand makes one think of the man who must have passed by. The smoke and the footprints are thus signs of something else. They are natural signs, inasmuch as the connection between the sign and the thing signified is a causal link, established by nature and learned from experience. These can be compared with road signs, for example, or such symbols as the outline of a heart pierced by an arrow. The connection between the symbol and the thing signified in these cases is not a natural one; it is established by human tradition or convention and is learned from these sources. These nonnatural signs, or symbols, are widely used in human communication.

In this framework the elements of language appear to be nonnatural signs. The interest in words and phrases reaches beyond their physical appearance: their perception is likely to direct attention or thought to something else. Words, in fact, are the chief media of human communication, and, as the diversity of languages clearly shows, the link involved between words and what they signify cannot be a natural one. Words and sentences are like symbols; they point beyond themselves; they mean something. Smoke means fire, the pierced heart means love. Words mean the thing they make us think of; the meaning of the word is the tie that connects it with that thing.

There are some words for which this approach seems to work very straightforwardly. The name Paris means (signifies, stands for, refers to, denotes) the city of Paris, the name Aristotle means that philosopher, and so forth. The initial plausibility of such examples created an obsession in the minds of many thinkers, beginning with Plato. Regarding proper names as words par excellence, they tried to extend the referential model of meaning to all of the other classes of words and phrases. Plato's theory of "forms" may be viewed as an attempt to find a referent for such common nouns as "dog" or for abstract nouns like "whiteness" or "justice." As the word Socrates in the sentence "Socrates is wise" refers to Socrates, for example, so the word wise refers to the form of wisdom. Unfortunately,

whereas Socrates was a real person in this world, the form of wisdom is not something to be encountered anywhere, at any time, in the world. The difficulty represented by "Platonic" entities of this kind increases as one tries to find appropriate referents for verbs, prepositions, connectives, and so forth. Discussion of abstract entities such as classes (*e.g.*, the class of all running things) and relations (*e.g.*, the relation of being greater than . . .) abound in philosophical literature; Gottlob Frege even postulated "the True" and "the False" as referents for complete propositions.

There are many more serious problems besetting the referential theory of meaning. The first one, eloquently pointed out by Frege, is that two expressions may have the same referent without having the same meaning. For example, "the Morning Star" and "the Evening Star" denote the same planet, yet, clearly, the two phrases do not have the same meaning. If they had, then the identity of the Morning Star and the Evening Star would be as obvious to anybody who understands these phrases as the identity of a vixen with a female fox or a bachelor with an unmarried man is obvious to speakers of English. As it is, the identity of the Morning Star with the Evening Star is a scientific and not a linguistic matter. Thus, even in the case of names, or expressions equivalent to names, one has to distinguish between the denotation (reference, extension) of the name—*i.e.*, the object (or group of objects) it refers to—and its connotation (sense, intension)—*i.e.*, its meaning.

The second problem with the theory of referential meaning arises from phrases that, though meaningful, pretend to refer but, in fact, do not. For example, in the case of such a definite description as "the present king of France," the phrase is meaningful although there is no such person. If the phrase were not meaningful, one would not even know that the phrase has no actual referent. Russell's analysis of these phrases, and the U.S. philosopher Willard V. Quine's similar treatment of such names as Cerberus, effectively detached meaning from reference by claiming that these expressions, when used in sentences, are equivalent to a set of existential propositions; *i.e.*, propositions without definite reference. For example, "The present king of France is bald" comes out as "There is at least, and at most, one person that rules over France, and whoever rules over France is bald." These propositions are meaningful, true or false, without definite reference.

Names, in fact, are very untypical words. The name of the third Secretary General of the United Nations, U Thant, has no meaning in English. Whether it means anything in Burmese does not matter either; the reference is not affected by the meaning or the lack of meaning of the name. Names, as such, do not belong to the vocabulary of a language; most dictionaries do not list them. Thus, in spite of the initial plausibility, the idea of reference does not help in understanding the nature of linguistic meaning.

Meaning and truth. Despite the failure of referential meaning, many philosophers were quite unwilling to give up the idea that the meaning of linguistic expressions has something to do with objects, events, and states of affairs in the world. They reasoned that if language is used to talk about the physical environment, then there must be some connection between man's words and the things around him. If reference fails to provide the link, something else must.

In the face of referential failure Russell fell back on truth. The Positivists suggested verifiability as the criterion of empirical meaning. Indeed, at least in many cases, it stands to reason that to understand a sentence is to know what state of affairs would make it true or false. Such considerations motivate the alethic (Greek *alētheia*, "truth") semantic theories, which claim that the notion of meaning is best explained in terms of truth rather than reference.

The most influential discussion of the notion of truth was offered by the Polish-born mathematician and logician Alfred Tarski in the 1930s. His semantic definition of truth is contained in the following formula (which he called [T]):

(T) X is true if, and only if, p

in which " p " is a variable representing any sentence and " X " is a variable representing the name, or unique de-

Influence
of transfor-
mational
grammar

Denotation
and con-
notation

Words as
symbols

Tarski's
semantic
definition
of truth

scription, of that sentence. The easiest way to obtain such a unique description is to put the sentence in quotation marks. Thus, we get such instances of (T) as

"Snow is white" is true if, and only if, snow is white.

The above formula implies a distinction between the object language and the metalanguage. "*X*" represents the name of a sentence in the object language—*i.e.*, roughly, the language used to talk about things in the world. The instances of (T) themselves, however, are in the metalanguage, the language in which one can talk about both things in the world and sentences of the object language. Tarski claimed that no language can contain its own truth predicate, for if it did then it would permit the formation of such sentences as:

(S) This very sentence is false.

Is the sentence (S) true or false? Clearly, it is true if, and only if, it is false, which is an intolerable paradox. Consequently, for any language the predicate "... is true" and other semantical predicates must belong to a language of a higher order (a metalanguage).

For this reason Tarski restricted his theory to clearly formalized, artificial languages, a decision that was very much in line with the positivistic tendencies of the 1930s. Nevertheless, the Tarski formula remained attractive even to some semanticists concerned with meaning in natural languages. For one thing, it seemed to succeed in pairing linguistic entities (named by the values of "*X*"; *e.g.*, the sentence "Snow is white") and nonlinguistic entities (named by the values of "*p*"; *e.g.*, the fact, or possible state of affairs, that snow is white). This correlation, however, is not very helpful because each one of the nearly infinite number of sentences one may form would have its "fact" as a counterpart, identifiable only by means of that very same sentence. Consequently, if linguistic meaning consisted in these correlations, no one could learn the meaning of any sentence at all and certainly not the meaning of all the sentences the speakers of a language are able to understand.

What was needed was a theory explaining the contribution of individual words—a clearly finite set—to the truth of sentences. Tarski himself, as well as other writers, suggested a repeatable procedure based on the notion of satisfaction. Snow, for example, satisfies the sentential function "*x* is white" because "Snow is white" is true. In much the same way, 3 satisfies the function " $2 \cdot x = 6$ " because " $2 \cdot 3 = 6$ " is true. Simply stated, the meaning of the predicate "... is white" is determined—and is learned—in terms of the set of objects of which it is true.

As this approach is extended to cover the wide variety of words that exist in a natural language, however, its initial simplicity—and thereby its attractiveness—becomes progressively lost. This can be illustrated by "egocentric" words like "I," "you," "here," and "now"; by connectives like "since," "however," and "nevertheless"; and, if these appear trivial, by such crucial words as "believe," "know," and "intend," on the one hand, and "good" and "beautiful" on the other. Whereas it is plausible to say that, for instance, Joe and Mary satisfy the function "*X* loves *Y*," provided Joe loves Mary, it is more complicated to determine what would satisfy such functions as "*X* believes *Y*," "*X* knows *Y*," or "*X* intends *Y*." "*X*" is satisfiable by people, but the satisfaction of "*Y*" poses a problem. If one suggests such things as propositions, facts, and possibilities, one is confronted with abstract entities of a kind similar to those encountered in the Platonic theory of referential meaning. Again, can it be said that John and a unicorn will satisfy the function "*X* looks for *Y*"—if it is true that John looks for a unicorn?

Another difficulty arises concerning "good," "beautiful," and other words of moral or aesthetic judgment. If, for example, beauty is indeed "in the eye of the beholder," then what one person calls beautiful might not similarly impress another, yet two people might keep arguing as to whether the thing is beautiful or is not. Thus, people may seem to agree on the meaning of the word, yet remain at odds about its proper application. The meaning of such "emotive" words cannot be decided in terms of truth alone.

A more serious objection to the alethic (truth) theory arises from the fact that many significant utterances of natural language are not true or false at all. Whereas statements, testimonies, and reports are true or false, orders, promises, laws, regulations, proposals, prayers, curses, and so forth are not assessed in terms of truth or falsity. It is not obvious that the employment of words in these speech acts is less relevant to their meaning than their use in speech acts of the truth-bearing kind.

Meaning and use. The difficulties just mentioned lead to another view concerning the notion of meaning, a theory that may be called the use theory. This view admits that not all words refer to something, and not all utterances are true or false. What is common to all words and all sentences, without exception, is that people use them in speech. Consequently, their meaning may be nothing more than the restrictions, rules, and regularities that govern their employment.

The use theory has several sources. First, in trying to understand the nature of moral and aesthetic discourse certain authors suggested that such words as "good" and "beautiful" have an emotive meaning instead of (or in addition to) the descriptive meaning other words have; in using them one expresses approval or commendation. If one says, for instance, that helping the poor is good, one does not describe that action, but says, in effect, something like "I approve of helping the poor, do so as well." Such is the role of these words, according to these thinkers, and to understand this role is to know their meaning.

The second, and more important, stimulus for the use theory was provided by the work of Ludwig Wittgenstein. This philosopher not only pointed out the wide variety of linguistic moves mentioned above but in order to show that none of these moves enjoys a privileged status proposed the idea of certain language games in which one or another of these moves plays a dominant or even an exclusive role. One can imagine, for instance, a tribe whose language consists of requests only. Members of the tribe make requests and the other members comply or refuse. There is no truth in this language, yet the words used to make requests would have meaning. Human language as it exists in reality is more complex; it is a combination of a great many language games. Yet the principle of meaning, according to this theory, is the same: the meaning of a word is the function of its employment in these games. To Wittgenstein the question "What is a word really?" is analogous to "What is a piece in chess?"

Finally, John L. Austin offered a systematic classification of the variety of speech acts. According to him, to say something is to do something, and what one does in saying something is typically indicated by a particular performative verb prefixing the "normal form" of the utterance. These verbs, such as "state," "declare," "judge," "order," "request," "promise," "warn," "apologize," "call," and so on, mark the illocutionary force of the utterance in question. If one says, for instance, "I shall be there," then, depending on the circumstances, this utterance may amount to a prediction, a promise, or a warning. Similarly, the words of the commanding officer, "You will retreat" may have the force of a simple forecast, or of an order. If the circumstances are not clear, the speaker always can be more explicit and use the normal form; *e.g.*, "I promise that I shall be there" or "I order you to retreat."

To rephrase the conclusion already stated: the dimension of truth and falsity is not invoked by all the utterances of the language; therefore, it cannot provide an exclusive source of meaning. There are other dimensions, such as feasibility (in case of orders and promises), utility (in case of regulations and prescriptions), and moral worth (in case of advices and laws). These dimensions may be as much involved in the understanding of what one said and, consequently, in the meaning of the words the speaker used, as the dimension of truth.

As previously mentioned, philosophers professing the alethic theory claimed that the meaning of a word should be explained in terms of its contribution to the truth or falsity of the sentences in which it can occur. The latest form of the use theory is an appropriate extension of the same idea. According to some exponents, the meaning of

Wittgenstein's language games

Limitations to the "meaning as truth" theory

Difficulties
with the
use theory

a word is nothing but its illocutionary act potential—*i.e.*, its contribution to the nature of the speech acts that can be performed by using that word. One difficulty with this view is that the definition is too broad, to the extent of being unilluminating or useless. Given this definition, nobody would know what any word means without knowing the entire language completely because the possibilities of employing a given word are not only without limit but extend to every conceivable context and circumstance. As Wittgenstein stated so forcefully,

The sign (the sentence) gets its significance from the system of signs, from the language to which it belongs. Roughly: understanding a sentence means understanding a language.

If this be the case, how can one account for the obviously gradual and prolonged process of learning a language? Indeed, the definition of the meaning of a word as illocutionary act potential seems to overstate the case. The obvious truth that the meaning of performative verbs, and other words closely tied to one illocutionary aspect or other, cannot be divorced from the nature of that type of speech act, does not entail that the meaning of an ordinary word, like "cat" or "running" is affected by any illocutionary force. Such words can occur in utterances bearing all kinds of illocutionary forces, so the contribution of these forces, as it were, cancel out. Nevertheless, what remains is the fact that all words are used to say something, in one way or another. The use theory would put a strong emphasis on the word "used" in the previous sentence. The next, and final, approach to meaning would stress the word "say."

Meaning and thought. In Wittgenstein's chess example, moves are made by moving the pieces. In a language, moves (saying something) are made by using words. And, according to the use theory, as a piece is defined by its move potential, so (the meaning of) a word is defined by its "saying" potential.

This analogy works only up to a certain limit. Whereas chess is only a game, the use of language is much more. One plays chess—or any other game—for its own sake; one speaks, however, with other ends in mind. Games, as it were, do not point beyond themselves; speech does. In order to see this, compare an ordinary conversation with such word games as children play—*e.g.*, exchanging words that rhyme, or words that begin with the same letter. These word games are language games and nothing more because the children use words according to certain rules. In doing so, however, they do not say anything except in the trivial sense of uttering words; nor are they called to understand what the other children say beyond the minimal feat of recognizing the words.

In real speech the situation is radically different. The point of using words in a real speech act is to be understood. If someone says, "It will rain tomorrow," his aim is to make the hearer believe that it will rain tomorrow. It is possible, of course, that the hearer or listener will not believe the speaker. Nevertheless, if the hearer understands what the speaker says, he will at least know that this is what the other person wants him to believe by saying what he says. Similarly, if one says, "Go home," and the listener understands what is said, then, whether or not the listener will actually go, he will at least know that this is what the speaker wants to bring about by using these words. Thus, the notion of saying something is inseparably tied to such concepts as belief, intention, knowledge, and understanding.

The view just outlined is a reformulation of a very traditional idea, namely, that speech is essentially the expression of thought. Words are used not to play a game with fixed rules but to express beliefs and judgments, intentions and desires; that is, to make others know, by the use of words according to fixed rules, that one has certain beliefs, desires, and so forth, and that one invites others to share them.

"Expression of thought" sounds rather vague. For one thing, what is a thought? Suppose John believes that Joe has stolen his watch. John can express this belief by saying, "Joe has stolen my watch" or "My watch has been stolen by Joe" or "It is Joe who has stolen my watch" and so on. Moreover, if John is a multilingual person, he can

express the same belief in German, French, and so forth. These variants, called paraphrases and translations respectively, will express the same belief, the same thought. But whereas it makes sense to ask for the exact words of John's statement or to ask about the language in which it was made, it would be foolish to ask for the exact words of John's belief or to ask about the language in which it is framed. The alternative in "Do you believe that Joe has stolen the watch, or that the watch was stolen by Joe?" does not make sense. Consequently, the same thought—the same proposition, as some philosophers prefer to call it—can be expressed by using various linguistic media. In other words, the same thought can be encoded in various codes (languages) and in various ways in the same code (paraphrases) in much the same way as the same idea can be expressed in speech or in writing and the same numbers can be written by using Roman or Arabic numerals.

From this point of view, it appears that saying something involves encoding a thought and that understanding what one said involves decoding and recovering the same thought. The meaning of a sentence will consist in its relation to the thought it is used to encode. This may be viewed as the fundamental thesis of the psychological theory of meaning.

As previously explained, no theory of meaning can be adequate as long as it treats sentences as indivisible units. For, in the first place, the potentially infinite number of sentences would defy any attempt to learn their meaning one by one, and, second, such a theory could not account for the obvious ability of fluent speakers to understand entirely novel sentences. There must be, therefore, a correlation between certain recurring elements of sentences (roughly, words) and certain recurring elements of thoughts (roughly, concepts or ideas). Accordingly, the learning of the semantic component of a language will consist in the learning of these connections.

In this learning process two notions play a prominent role: synonymy and analyticity. As the sentences that express the same thought stand in the relation of paraphrase (or translation), so words or phrases that code the same idea are related as synonyms—*e.g.*, "vixen" and "female fox." Again, because one concept may include another, the sentence expressing this relation will record a conceptual truth or analytic proposition; *e.g.*, "A dog is an animal." A definition, finally, will exhibit all parts of a concept by a combination of such propositions.

What concepts are and how they are related to words are topics that have been discussed throughout the history of philosophy. The following problems related to concepts pertain to the core of philosophical psychology: whether all concepts are derived from experience, as Aristotle and the Empiricists believed, or whether some of them at least are innate, as Plato and the Rationalists maintained; whether concepts exist prior to and independent from their verbal encodings, as the Realists and Conceptualists claimed, or whether they are nothing but a certain "field of force" accompanying the words, as the Nominalists thought.

It should be noted that these disputes, in a modern garb, still continue with undiminished force. Contemporary Empiricists still try to reduce most concepts to a configuration of sense data, or a pattern of nerve stimulation, while the Behaviourists attempt to explain understanding in terms of overt behaviour. Modern-day Rationalists reply by insisting on the unique spontaneity of human speech and by reviving the theory of innate ideas.

MEANING IN LINGUISTICS

Semantics in the theory of language. The science of linguistics is concerned with the theory of language expressed in terms of linguistic universals—*i.e.*, features that are common to all natural languages. According to the widely adopted schema of the U.S. scholar Charles W. Morris, this theory must embrace three domains: pragmatics, the study of the language user as such; semantics, the study of the elements of a language from the point of view of meaning; and syntax, the study of the formal interrelations that exist between the elements of a language (*i.e.*, sounds, words) themselves. Subsequently, certain authors spoke of three levels: the phonetic, the syntactic (the pho-

Synonyms
and
analytic
proposi-
tions

Speech
as an
expression
of thought

Prag-
matics,
semantics,
and syntax

netic and syntactic together are often called grammatical), and the semantic level. On each of these levels a language may be studied in isolation or in comparison with other languages. In another dimension, the investigation might be restricted to the state of a language (or languages) at a given time (synchronic study), or it might be concerned with the development of a language (or languages) through a period of time (diachronic study).

Semantics, then, is one of the main fields of linguistic science. Yet, except for borderline investigations, the linguist's interest in semantic matters is quite distinct from the philosopher's concern. Whereas the philosopher asks the question "What is meaning?", the typical questions the linguist is likely to ask include: "How is the meaning of words encoded in a language?" "How is this meaning to be determined?" "What are the laws governing change of meaning?" and "How can the meaning of a word be given, expressed, or defined?"

A few examples will suffice to illustrate some of these problems, and to show how the linguist's approach differs from that of the philosopher. In the matter of encoding, words are arbitrary signs; to some authors, particularly to the Swiss linguist Ferdinand de Saussure, this feature of arbitrariness represents an essential characteristic of all real languages. Nevertheless, in all languages there are clear cases of onomatopoeia—*i.e.*, the occurrence of imitative words, such as "whisper," "snore," "slap," and, more remotely, "cuckoo."

There are several other issues that pertain to the question of encoding. Certain languages show a marked preference for very specific words, at least in certain domains, while lacking the corresponding general terms, which are the only ones occurring in other languages. The Eskimos, for instance, have a number of words denoting various kinds of snow, but no single word for snow. Similarly, in English, although there are distinct names for hundreds of animal species, there is no name for the very familiar animal species of which the female member is called cow and the male member bull.

There are also languages, such as English or Chinese, that for the most part prefer single words to the compounded phrases that other languages (*e.g.*, German) seem to favour. Accordingly, whereas the English vocabulary is larger, the German words are more pliable, capable of entering into compounds often of great length and complexity. Such differences support Saussure's distinction between lexicological and grammatical languages.

Another distinction can be drawn concerning the relative frequency and importance of context-bound and context-free words. The meaning of such English words as "take," "put," and "get" depends almost entirely on the context—*e.g.*, putting up with somebody has very little to do with putting off something or other. These can be compared with verbs like "canter" or "promulgate," which, by their very specific meaning, almost determine the context, rather than having their meaning determined by the context. Clearly, the context-bound type of word, such as "take" or "put," lends itself to idiomatic use, rather than the context-free word.

There are some obvious regularities in the change of meaning that are of interest to the linguist. One such regularity is the extension or transference of meaning based upon some similarities—*i.e.*, the phenomenon of metaphor. For example, one can speak of the leg of the table, the mouth of the river, the eye of the needle, and the crown of the tree. These are anthropomorphic metaphors: the transfer goes from something belonging to an individual or close to him (his body, garments) to something more remote. The same principle operates in the extension of meaning from a domain close in interest rather than in physical proximity. Baseball-minded people are apt to speak of "not getting to first base," "striking out," or "scoring a hit" in contexts often remote from baseball. For a similar reason, many abstract concepts are denoted by words transplanted from the concrete domain. Such phrases as "grasping ideas," "seeing the point of a joke," "body of knowledge," "in the back of my mind," and many others, are the result of this very important move from the abstract to the concrete.

Meaning changes of another type are the result of emotive factors. The word democracy, for instance, has all but lost its original meaning and has become a word applicable to any system the speaker wants to praise. The contrary development is exhibited in the recent history of such words as "Fascist" and "aggression." In order to avoid derogatory connotations one is often forced, by social pressure, to use euphemisms, often to the detriment of clarity. Examples of this include the switch from "underdeveloped nations" to "developing nations," from "retarded children" to "exceptional children," and from "old people" to "senior citizens."

The preceding are but a few examples concerning coding and meaning change. The questions about the ways of finding out what a word means and about the manner of giving an adequate definition of a word deserve a more detailed account.

Meaning, structure, and context. Foreigners in a strange country and linguists are often confronted with the task of learning a new language. It is important to realize that in doing so they do not set out with a completely blank mind: they expect to learn a language (*i.e.*, a system of communication describable in terms of a large set of linguistic universals). They expect—consciously in the case of the linguist and unconsciously in the case of the layman—to find words and sentences, grammatical structures, and illocutionary forces in that language. And, on the semantic level, they expect words that will fit into the familiar semantic classes. They are confident, in other words, that the language they intend to master will be intertranslatable with their own.

Therefore, although at the very beginning their learning remains on the ostensive level (trying to find out the name of this or that kind of object), very soon they proceed to the level of first guessing, then establishing, the meaning of words from the contexts in which they occur. This has to be the case, for words that in any way can be viewed as "names" of objects (*i.e.*, that could be learned ostensively) form but a fraction of the vocabulary of any language. Anyone who doubts this should but try to list the words from this paragraph that could be learned ostensively. Moreover, linguists find no great difficulties in learning dead languages—*e.g.*, that of the ancient Egyptians—without any contact with any speaker, provided that a sufficiently large corpus of texts is available and that some clues are provided to the meaning of at least some words.

If any more evidence concerning this point is needed, one should remember that "pictorial" dictionaries are bound to remain on the kindergarten level, and that the mark of a good dictionary is the abundance of appropriate contexts. Thus, the contexts show the concept.

These intuitions are behind the U.S. philosopher Paul Ziff's semantic theory. According to Ziff, the meaning of a word is a function, first, of its complementary set, which consists of all the acceptable sentences in which the word can occur, and, second, of its contrastive set, which consists of all of the words that can replace that word in all of these sentences without rendering the sentences deviant. Clearly, the elaboration of the contrastive set will produce words more and more similar in meaning to the word in question, the limiting case being synonyms that can occur wherever the word in question can occur.

This theory is in need of further refinement. In the sentence "The cat sleeps," the fact that "cat" can co-occur with "sleep" undoubtedly casts some light on the meaning of these words (a cat is a kind of thing that can sleep). But there are a great number of sequences that could complete the frames "The cat sleeps and . . ." ". . . said that the cat sleeps," and so forth. Clearly, the near infinity of the resulting sentences will not contribute anything to the meaning of "cat" beyond what the segment "the cat sleeps" already contributes.

Transformational grammar can be of assistance at this point. According to this approach, the sentences just considered are simply surface forms, each corresponding to an underlying structure, in which "cat" and "sleep" appear as forming an elementary, or kernel sentence (roughly: "a cat sleeps"). The essence of Ziff's insight can be rein-

Structure and language learning

Context-bound and context-free words

Surface forms and underlying structures

terpreted in terms of the notions developed by Zellig S. Harris: co-occurrence (instead of complementary set) and co-occurrence difference (instead of contrastive set), both restricted to kernels. Because the vocabulary of a language is limited and the number of kernel structures is very small, the meaning of a word can be determined on the basis of a finite set of elementary sentences.

The contribution of grammar to semantic theory is by no means exhausted by this step. For the grammatical restrictions on a word represent, as it were, the "skeleton" of its meaning before the "flesh" is put on by the co-occurrences. The very first step in giving the meaning of a word is to specify its grammatical category—noun, verb, adjective, adverb, connective, and so forth—and not to speak of grammatical constants (such as the first, but not the second, "to" in "I want to go to Paris"), the meaning of which, if any, is entirely determined by their grammatical role. A refined grammar yields much more: the fact that the adjective "good," for example, unlike adjectives like "yellow" or "fat," can occur in the frames "(He is) good at (playing chess)"; "(The root is) good to (eat)"; "It is good that (it is raining)"; "It was good of (him) to (come)" says a great deal about the meaning of that word. The co-occurrences then complete the picture.

Lexical entries. Good dictionaries offer a variety of contexts for the items listed, but, obviously, this is not enough. For one thing, no dictionary can list all the co-occurrences. There must be devices to sum up, as it were, the information revealed by the contexts. This is the role of dictionary definitions. The branch of scientific semantics that is concerned with the form and adequacy conditions of dictionary entries is called lexicography.

A systematic study of dictionary entries was presented in the 1960s by the U.S. philosophers Jerrold J. Katz and Jerry A. Fodor. According to them, the standard form of a dictionary entry comprises three kinds of ingredients: grammatical markers, semantic markers, and distinguishers. The grammatical markers describe the syntactic behaviour of the item in question in terms of a refined system of grammatical categories. The traditional division of words into nouns, adjectives, verbs, adverbs, and so on is but the first step in this direction. The class of nouns, for example, has to be subdivided into count nouns (like "cat"), mass nouns (like "water"), abstract nouns (like "love"), and so forth. The class of adjectives must be classified into subclasses that are fine enough to capture the grammatical peculiarities of such adjectives as "unlikely" or "good." The traditional subdivision of verbs into transitives and intransitives has to be completed to account for such verbs as "compare" or "order," which obviously involve three noun phrases ("someone compares something to something"), often of a particular kind ("human" nouns or noun clauses).

The idea of a semantic marker is merely a further elaboration of the traditional notions of genus and species. The result is a system of semantic markers that comprise such items as "physical object," "animate," "human," "male," "young" (in the case of the entry for "boy"), and others. Katz claims that the problems of synonymy, analyticity, and contradiction can be handled, at least in part, in terms of lexical items sharing some or all of their semantic markers.

Finally, the distinguisher completes the dictionary entry by giving, as it were, the leftover, if any, of the semantic information. There is no general form for the distinguisher; it may give the atomic weight (for elements), purpose (for tools), concise description (for animals), and so forth.

Generative semantics. According to the original formulation of generative or transformational grammar, the semantic and the syntactic components were regarded as distinct elements in the deep structure of a sentence. The syntactic component consisted of a relation of phrase markers giving the transformational structure of the sentence, usually represented in terms of "tree" diagrams with such nodes as "S" for sentence, "NP" for noun phrase. The semantic content entered through the process of lexical insertion—*i.e.*, the replacement of some of the nodes by words, the carriers of meaning. Lexical insertion was

supposed to take place at the very beginning of the series of transformations leading up to the surface form of the sentence. The original input of meaning, as it were, was carried through the transformations yielding the semantical reading, or sense, of the whole sentence.

Several modern studies have attempted to demonstrate that this separation between syntax and semantics cannot be maintained. It appears that certain words in themselves indicate a structure analogous to syntactic structures. For example, consider "harden" and "break." To harden something is to cause that thing to become hard (or harder); to break something is to cause something to become broken. Because "harden" consists of two elements, "hard-" and "-en," thus it could be argued that the word itself is structured; "break" does not indicate any structure, yet its meaning clearly involves one. "Broken," therefore, carries a more basic semantic unit than "break." Again, in the case of such verbs as "remind," "allege," "blame," or "forgive," one feels that their meaning is highly structured, and, with some thinking, one can articulate the presuppositions of these words, which involve a great number of semantic units in a very complex relationship.

These and similar reasons support the theory of generative semantics, which denies a clear distinction between the semantic and the syntactic components. The transformations, in this theory, connect the surface structure of the sentence with its semantic representation (or according to some linguists, its logical form). Words, then, can encode either a semantic primitive (such as "blue") or a whole structure (such as "forgive") within the semantic representation. Thus, there is no definite point in the transformational history of the sentence at which lexical insertion must occur. The ultimate conclusion of this view is that, instead of the threefold division of semantics, syntax, and phonetics, all that is needed is the simple distinction between semantics and phonetics, corresponding to the distinction between meaning (as structured) and its verbal encoding. How much, finally, of the semantic structure can be attributed to a particular language, and how much can be ascribed to common (and possibly innate) elements of the human mind, remains a fascinating problem for continued study. (Z.V.)

The study of writing

This section deals with the various fields of study relating to written sources and writing systems. Of these, the most important are philology, which deals mainly with written and oral source materials; and linguistics and grammatology, which deal mainly with systems of signs, which in turn are based on primary source materials. The use of the word "mainly" implies that, while certain fields deal predominantly with certain aspects of study, there are no sharp divisions between them, and some areas of study can be treated by both disciplines.

THE SCOPE OF THE STUDY OF WRITING

Sources of information. In a study of any system of writing, a basic distinction must be made between the raw materials (subject matter) that are studied and the systems deduced from them. Spoken (or written) utterances represent the raw materials of a language; a grammar or a lexicon presents a linguistic system abstracted and reconstructed from these materials. Similarly, written texts represent the raw materials of writing from which an alphabet can be reconstructed.

Manifold sources of information bear on systems of writing: (1) Written sources proper include texts, inscriptions, books, and manuscripts. The ability to understand these sources can be handed down traditionally from generation to generation, as in the case of the Hebrew or Chinese writings, or it must be recovered by a process of decipherment, as in the case of the Egyptian hieroglyphic. (2) Lists of signs, alphabets, syllabaries, and so on can be devised by the "inventors" of a writing at the time when the writing is first introduced; or they can be reconstructed in later times from the written sources by teachers or scholars for didactic or scholarly purposes. (3) Studies are based on

Written sources

written sources or lists of signs, alphabets, syllabaries, and so on. Such studies can be primary, when they are based directly on sources, or secondary, when they must rely on primary studies for information.

A partial study, which is delimited by the subject matter, times, and area, must be distinguished from a general, comprehensive study of the subject. Thus an analysis of the writing of Middle English, Egyptian hieroglyphic writing, or the Greek alphabet may be contrasted with an analysis of the structure of alphabetic writings or a study of cursive or monumental writing.

Definitions of terms. *Semiotics.* As has been noted above, men communicate with each other by means of various systems of signs, of which the most universal are language, a system of auditory communication; and gesture language and writing, two systems of visual communication. For the general science of signs, several terms have been proposed, of which the term semiotic, or semiotics, as here preferred, may perhaps be the most appropriate. Semiotics covers a much broader area than the term semantics, which deals with the meaning of linguistic elements.

Philology, linguistics, and grammarology. Philology, involved mainly in the study of the linguistic sources of a people or a group of peoples, forms the basic means for the comprehension of their respective cultures. It deals less with oral sources than with written sources, mainly literature (whatever its exact meaning). Philology deals with the formal aspects of writing under the topic of epigraphy and paleography. Linguistics is concerned with the study of linguistic systems as reconstructed mainly from oral sources. Pursued less than the study of "oral language," the study of the "written language"—that is, of the language as it is used in written sources—is also a matter of linguistics. Linguistics deals with the structural aspects of writing under the heading of graphemics.

The field of study that deals with writing in the broadest sense is called grammarology. Equally appropriate terms for this subject are grammatonomy and graphonomy.

GRAMMATOLOGY

Three main approaches. Three main approaches to grammarology can be distinguished: descriptive-historical, typological-structural, and formal.

Descriptive-historical. The traditional descriptive-historical approach to the study of writing is by far the most common. This is a simple narrative approach to the description of writing in its historical evolution. The apparent shortcoming of descriptive-historical texts is the general lack of systematic typology—that is, systematic classification by type. Good studies on individual writings, such as hieroglyphic Egyptian and the Greek alphabet, are not wanting. What is entirely missing is theoretical and comparative evaluation of the various types of writing, such as discussions of various types of syllabaries, alphabets, word signs, and logo-syllabic writings.

The historical approach is further vitiated through confusion with considerations of a geographic nature, as evidenced by such chapter titles as "Asiatic Writings" and "American Writings," or "Writings of Asia" and "Writings of America," which are frequently found in the standard manuals on writing.

Typological-structural. The typological-structural approach is based on the realization of the importance of structure and typology in the study of writing. In contrast to the traditional approach, according to which the writings of the world are described in their evolutionary progress, the new approach requires first a thorough analysis of the structure of the individual writing systems and then their classification by type within the framework of writing in general.

Under "structure" is meant here primarily what is sometimes called the "inner structure," which is concerned with the function of a writing system, in contrast to the "outer structure," which involves its formal characteristics. Considerations of function are based on such points as word writing (logography), syllable writing (syllabography), and letter writing (alphabetography), and the typology of logo-graphic, syllabic, and alphabetic systems of signs. Con-

siderations of form involve such points as the contrast between the pictorial and linear systems or between the monumental and cursive writings. While, theoretically, both approaches are acceptable, the emphasis placed here on function, rather than form, may be illustrated by considering the following case. From the point of view of function, the Morse alphabet represents the same type of alphabet as the Latin writing and its descendants: from the point of view of form, the Morse alphabet is formally independent of the Latin writing. Based on considerations of function, the Morse alphabet is considered as being of the Latin type, despite its different formal, outer structure.

Formal. There are two kinds of formal approach to the field of grammarology: the traditional approach, as practiced mainly in the philological disciplines under the topic of epigraphy and paleography (see below), and the formal approach to sign analysis recently initiated in the United States. The aim of the latter is to provide a scientific account of cursive writing as practiced in that country. Its procedure consists first of breaking up the letters of the cursive writing into its component segments, which appear in the form of bars, hooks, arches, and loops, and then of providing formation rules governing the linking together of these segments into letters and into larger strings corresponding to words of the language. It is said that 18 such segments are sufficient to describe all English lower- and upper-case letters.

Subdivisions of grammarology. Three main subdivisions of grammarology can be distinguished: subgraphemics, graphemics, and metagraphemics. They are treated mainly by scholars versed in philological and linguistic disciplines.

Subgraphemics. The field of subgraphemics deals mainly with primitive forerunners of writing that utilize visual marks having no set correspondences in language. In its sublinguistic aspects, subgraphemics can be compared to kinesics, the study of the various communicational aspects of learned, patterned body motion behaviour; and cherology, the study of gesture language.

Graphemics. The field of graphemics deals with full writing or phonography, as represented in systems of writing in which written signs generally have set correspondences in elements of language. The field of graphemics thus deals with writing after it became a secondary transfer of the language, a vehicle by which elements of the spoken language were expressed in a more or less exact form by means of visual signs used conventionally. This took place for the first time about 5,000 years ago in the Sumerian and Egyptian writings.

Instead of "graphemics," other scholars use the terms "graphics" or "graphic linguistics." All three terms are frequently misused by scholars who limit the terms to the study of alphabetic writings, overlooking or paying scant attention to all other types of writing, such as the logosyllabic and syllabic systems.

Little work has been done in the field of relations of writing to language. Philologists have been concerned mainly with the historical evolution of writing and have paid little attention to the interrelations between writing and language. Linguists have been more concerned with the spoken language than with the written language. When interested in written languages, they have often limited their study to living written languages, neglecting the rich sources of information that can be culled from ancient written languages and from pre-alphabetic systems. The question of the relationship of writing to language has been pursued in recent years mainly by scholars with a background in linguistics. Because of their interest in modern languages and writings, this implies generally relations between the alphabet and language. A general treatment of the subject can be found in the respective chapters of the introductory manuals to linguistics. Linguists generally have stressed the independent character of writing and have studied it as an independent system rather than as a system ultimately based on and related to the underlying language.

While the connections between language and writing are close, there has never been a one-to-one correspondence between the elements of language and the signs of writing. The "fit" (*i.e.*, the correspondence) between language and

The growing divergence of language from writing

writing is generally stronger in the earlier stages of a certain system of writing and weaker in its later stages. This is because a writing system when first introduced generally reproduces rather faithfully the underlying phonemic structure (structure of sounds). In the course of time, writing, more conservative than language, generally fails to keep up with the continuous changes of language and, as time progresses, diverges more and more from its linguistic counterpart. A good example is the old Latin writing, with its relatively good "fit" between graphemes (the written letters or group of letters that represent one phoneme or individual sound) and phonemes as compared with the present-day French or English writing, with their tremendous divergences between graphemes and phonemes. In some cases, recent spelling reforms have helped to remedy the existing discrepancies between writing and language. The best "fit" between phonemes and graphemes has been achieved in the Korean writing in the 16th century and in the Finnish and Czech writings of modern times.

Families of writings are not related to families of languages. Note, for example, that English and Finnish are written in the Latin writing but belong to two different families of languages, and that the cuneiform writing was used in antiquity by peoples speaking many different languages.

The temporal primacy of language over writing has been taken for granted by most scholars, especially the American linguists. It has been contested by some European scholars, who claim that writing is as old as oral language and gesture language. The fact is that full writing, expressing linguistic elements, originated only about 5,000 years ago in Mesopotamia and Egypt and that full writing is therefore much younger than language. Only if the semasiographic stage is included under writing can the assumption of equal temporal hierarchy of writing and language be admitted. (Semasiography is the use of marks to convey meaning without the presence of linguistic elements.) As noted elsewhere, however, the semasiographic stage should not be treated as full writing, but as a forerunner of writing.

Metagraphemics. A study of the various metagraphic devices (e.g., punctuation marks, capital or italic letter forms) that are used besides or in addition to writing proper may be called metagraphemics or paragraphemics. This is still an obscure field, and its relationship to both subgraphemics and graphemics needs a thorough investigation.

EPIGRAPHY AND PALEOGRAPHY

The investigation of writing from the formal point of view has been traditionally the prime domain of the epigrapher and paleographer. Epigraphy is concerned mainly with inscriptions written in characters that are incised or scratched with a sharp tool on hard material, such as stone or metal; paleography deals mainly with manuscripts written in characters that are drawn or painted with pen, pencil, or brush on soft material, such as leather, papyrus, or paper. Since epigraphy means "writing upon something" and paleography means "old writing," it is clear that the distinction made above between epigraphy and paleography cannot be justified on etymological grounds. The distinction has grown artificially over the years, as one scholar or another began to apply one or the other term to his own branch of study of written sources. Because of the close interrelations between epigraphy and paleography, some scholars refuse to admit any distinction between the two and prefer to use only the term paleography.

The main characteristics of epigraphy and paleography as listed above may be applied, with some leeway, to the ancient Near East (Mesopotamia, Egypt, Anatolia), the classical world, China, India, the Islamic world, and, in general, to the Western writings from the Middle Ages down to the introduction of the printing press. But there are some difficulties: Mesopotamian and Aegean clay tablets are soft, and the writing is cursively executed, both points characteristic of paleography, but the tablets are also durable and bulky and have incised, concave characters executed with a stylus, all points characteristic of epigraphy. Similarly, the wax tablets of the classical world

are soft, perishable, and cursively executed, characteristics of paleography, but are written with concave characters, incised with a stylus, characteristic of epigraphy. There are likewise some difficulties in the classification of tablets of wood; they are soft and perishable, and the writing on them is generally cursive, characteristic of paleography, but they have characters incised or scratched with a sharp implement, characteristic of epigraphy.

Paleography and epigraphy are involved in the study of written sources from two points of view: the purely formal aspect and the hermeneutical (interpretive) aspect.

The study of the purely formal aspect, possible without any understanding of the contents or without an extended study of the contents, is concerned, for example, with the kind, form, and size of the materials; the technique of writing; and the form, order, and direction of writing. Hermeneutics, possible only with study of the contents, is concerned, for example, with the dating and localizing of written sources, their authorship, linguistic interpretation, and content evaluation.

A general scientific discipline of epigraphy and paleography does not exist. There are no studies that treat of the subject from a general, theoretical point of view, encompassing all the written sources, wherever they may be found. There are, for example, no treatises listing and discussing the various materials, or shapes and sizes of materials, used for writing throughout the world, just as there are no structural-typological studies that treat of the formal evolution of signs from pictorial to linear or from round to angular. Among other potential topics that await investigation are: trends in ductus (hand; the general shape and style of letters), such as individual, national, and regional; the direction of writing, indication of prosodic features (quantity, stress, and tone), and names of signs (letters). The narrow fields that are represented are, for example, West Semitic epigraphy, Arabic paleography, Greek and Latin epigraphy and paleography, or Chinese epigraphy and paleography. In all cases, these narrow fields of study form subdivisions of wider but still linguistically or geographically defined fields of study, such as Semitic or Arabic philology, classical philology, Assyriology, and Sinology.

HISTORY OF THE STUDY OF WRITING

The first students of writing were doubtless the very originators ("inventors") of a new writing system. By "writing system" is meant here a full writing in which the individual signs of the writing stand for the corresponding elements of the language, which is to be contrasted with the forerunners of writing, in which the individual signs have but loose connection with language. As a result of a discrete analysis of the language for which a writing system was devised, lists of the elements of a language and their proposed written counterparts must have been first compiled and experimented with in actual practice. The establishment of a full system of writing also required conventionalization of forms and principles. Forms of signs had to be standardized so that the users would draw the signs in approximately the same way. Regulation of the system had to take place in the matter of the orientation of signs and the direction, form, and order of the lines, columns, and the sides of a text. Correspondences established between signs and words were paralleled by those between signs and definite syllabic values. After the initial period of trial and error, the established correspondences were conventionalized by being taught in schools.

Studies prior to the 18th century. The activities involved in setting up a full system of writing are indirectly attested in the Sumerian school texts, known almost from the beginnings of the Sumerian writing, which appear in the form of lists of signs and words, and scribal and literary exercises. The scribal activities of the Sumerians and, in the later periods, Akkadians are matched, albeit to a smaller degree, as far as actual attestation is concerned, by those of other peoples of the ancient Near East, such as the Egyptians and Hittites.

The study of the language and the corresponding writing was highly developed among the Chinese and Indic peoples, as best exemplified by the great Indic grammarian Pāṇini (about the 4th century BC) and his school, as

Formal and interpretive studies

Definitions of epigraphy and paleography

The Sumerian school texts

well as among the Greeks and to a lesser degree among the Romans. Beginning in the Middle Ages, the Arabs and Jews showed great interest in matters pertaining to their languages and writings. Important contributions were made by the Arab scholars Sibawayh (8th century) and az-Zamakhshari (1075–1143) and by the Jews Rashi (1040–1105) and David Kimhi (c. 1160–c. 1235). Interesting, but largely fantastic, is a collection of several dozen alphabets and sign lists put together by the Arab Ahmad ibn Abū Bakr ibn Waḥshīyah (c. 800). Among the early Europeans were the Spanish bishop Diego de Landa (1524–79), with his analysis of the Maya writing, and the German scholar Athanasius Kircher (1602–80), with his frequently mystic ideas about writing, especially the Egyptian hieroglyphic.

Modern Western studies. Modern general studies of writing in the West began in the second half of the 18th century. This group of early studies was based almost exclusively on Greco-Latin writing, with its further developments in the Middle Ages and modern times, and supplemented by a scattering of Semitic alphabets, such as the Hebrew, Arabic, and Syriac. The books of this first period are of no more than historical interest today.

In contrast to these early studies, the second group consists of the larger and much more serious undertakings of François Lenormant (1872), Heinrich Wuttke (1872), Carl Faulmann (1880), Isaac Taylor (1883), and Philippe Berger (1891). These books became standard manuals in the field of writing and served that function until the first half of the 20th century.

Next in time is a group of studies generally smaller in size and more limited in coverage than the group discussed just above, as represented by Walter James Hoffmann (1895), R. Stübe (1907), Theodore Wilhelm Danzel (1912), Karl Weule (1915), and William A. Mason (1920). On the basis of such monumental works on the American Indian writings as those of Henry R. Schoolcraft (1851), Garrick Mallory (1886 and 1893), and other similar works, the authors of this group of studies have emphasized much more than did the earlier scholars the importance of the forerunners of writing in the whole field of the study of writing.

The main characteristic of the next period, the first half of the 20th century, during which great manuals on writing were produced, was the descriptive-historical approach, as it was of all the preceding group of studies on writing.

The typological-structural approach to the study of writing appeared in the latter half of the 20th century, along with interest in the relationship of writing to society. Representatives of this and of the previous group are briefly characterized in the bibliography. (I.J.G./Ed.)

BIBLIOGRAPHY

General linguistics: ROBERT H. ROBINS, *A Short History of Linguistics*, 2nd ed. (1979), *General Linguistics: An Introductory Survey*, 3rd ed. (1980), is a comprehensive and balanced treatment of the whole field. LEONARD BLOOMFIELD, *Language* (1933), a classic introduction to the subject, is still not completely superseded and is essential reading for an understanding of subsequent American work. CHARLES F. HOCKETT, *A Course in Modern Linguistics* (1958), a comprehensive, stimulating, though somewhat personal textbook, represents the post-Bloomfieldian period in the United States. JOHN LYONS has

produced a number of notable surveys: *Introduction to Theoretical Linguistics* attempts to synthesize more traditional and more modern ideas on language, paying particular attention to generative grammar and semantics; *New Horizons in Linguistics* (ed., 1970), contains previously unpublished chapters on developments in most areas of linguistics; *Language and Linguistics: An Introduction* (1981) is a textbook covering theoretical developments. MARTIN JOOS (ed.), *Readings in Linguistics* (1957), is an excellent selection of key articles on structuralism in the post-Bloomfieldian period. Z.S. HARRIS, *Methods in Structural Linguistics* (1951), offers the most extreme and most consistent expression of the distributional approach to linguistic analysis—important for the development of generative grammar. NOAM CHOMSKY, *Syntactic Structures* (1957), is the first generally accessible and relatively non-technical treatment of generative grammar, widely recognized as one of the most revolutionary books on language to appear in the 20th century; J.P.B. ALLEN and PAUL VAN BUREN (eds.), *Chomsky: Selected Readings* (1971), contains an annotated selection of key passages from Chomsky's main works. S. PIT CORDER (ed.), *The Edinburgh Course in Applied Linguistics*, 4 vol. (1973–77), is a collection of readings covering a wide range of views. RICHARD C. OLDFIELD and J.C. MARSHALL (eds.), *Language* (1968), J.A. FODOR, T.G. BEVER, and M.F. GARRETT, *The Psychology of Language* (1974); and JOSEPH F. KESS, *Psycholinguistics* (1976), are important works in psycholinguistics. DELL HYMES (ed.), *Language in Culture and Society* (1964), is an excellent selection of articles in sociolinguistics and anthropological linguistics. Valuable information is found in ROY HARRIS, *The Language Myth* (1981); GEORGE A. MILLER, *Language and Speech* (1981); SUSAN BASSNETT-MCGUIRE, *Translation Studies* (1981); ERIC WANNER and LILA R. GLEITMAN (eds.), *Language Acquisition: The State of the Art* (1982); and HANS AARSLEFF, *From Locke to Saussure: Essays on the Study of Language and Intellectual History* (1982).

Semantics: W.P. ALSTON, *Philosophy of Language* (1964), is the best current introduction to philosophical semantics. STEPHEN ULLMANN, *Semantics: An Introduction to the Science of Meaning* (1962), has become a classic work. M. BLACK, *Language and Philosophy* (1949), discusses some earlier views. L. BLOOMFIELD, *Language* (1933), contains a classic discussion of scientific semantics. B.L. WHORF, *Language, Thought and Reality*, ed. by J.B. CARROLL (1956), raises the issue of linguistic relativism. J.J. KATZ, *The Philosophy of Language* (1966), offers a semantic theory tied to generative grammar, the best expression of which is found in N. CHOMSKY, *Aspects of the Theory of Syntax* (1965). W.V. QUINE, *Word and Object* (1960); and P. ZIFF, *Semantic Analysis* (1960), represent two different but influential semantic theories. Philosophy of language is explored also in the following monographs: IRMENGARD RAUCH and GERALD F. CARR (eds.), *The Signifying Animal: The Grammar of Language and Experience* (1980); CHARLES ALTIERI, *Act and Quality: A Theory of Literary Meaning and Humanistic Understanding* (1981); GRAHAM D. MARTIN, *The Architecture of Experience: A Discussion of the Role of Language and Literature in the Construction of the World* (1981); NATHAN U. SALMON, *Reference and Essence* (1981); DERECK BICKERTON, *Roots of Language* (1981); SÁNDOR HERVEY, *Semiotic Perspectives* (1982); JEREMY CAMPBELL, *Grammatical Man: Information, Entropy, Language, and Life* (1982); ANNETTE LAVERS, *Roland Barthes, Structuralism and After* (1982); and DAVID LIGHTFOOT, *The Language Lottery: Toward a Biology of Grammars* (1982).

Periodicals: *Language, Word, International Journal of American Linguistics* (United States); *Philological Society Transactions, Journal of Linguistics* (Great Britain); *Lingua. Studies in Language* (Holland); *Bulletin de la Société de Linguistique de Paris* (France).

(Ed.)

Lisbon

Lisbon (Portuguese: Lisboa), the capital of Portugal and of the district that bears its name, is the nation's chief port, largest city, and commercial, political, and tourist centre. It stands on the westernmost point of land of continental Europe. The city's name is a modification of the ancient Olisipo (variant Ulyssipo), and its founding has been variously attributed to Ulysses (Greek: Odysseus), the hero of Homer's *Odyssey*, to Elisha, grandson of the Hebrew patriarch Abraham, and, more credibly, to Phoenician colonists.

Lisbon owes its historical prominence to its superb natural harbour, one of the most beautiful in the world. The city lies on the north bank of the Tagus River (Rio Tejo), about eight miles (13 kilometres) from the river's entrance into the Atlantic Ocean. From the ocean upstream to the city, the river is almost straight and about two miles wide. It is spanned, on the west side of the city, by the 25th of April Bridge (formerly called the Salazar Bridge), the longest suspension bridge in western Europe. Just east of the bridge, the Tagus suddenly broadens into a bay seven miles wide called the Sea of Straw (Mar de Palha)—a reference to the sheen of the water. Scenically spectacular though it may be, this hill-cradled bay of burnished water lies on a strategic sea route and serves as a busy port, handling much of the exports and imports of Portugal and Spain.

This article is divided into the following sections:

Physical and human geography	72
Character of the city	72
The landscape	72
The city site	
Climate	
The city layout	
Housing	
The people	73
The economy	73
Industry	
Commerce and finance	
Transportation	
Administration and social conditions	74
Government	
Health	
Education	
Cultural life	74
History	74
The early period	74
Prehistoric to Moorish times	
The Portuguese conquest	
The Age of Discovery	
Evolution of the modern city	75
Disaster and reconstruction	
19th-century expansion	
The 20th century	
Bibliography	76

Physical and human geography

CHARACTER OF THE CITY

Once a remote outpost on what was thought to be the farthest edge of the known world, by the 15th century Lisbon was established as the centre of operations for Portuguese exploration. Although this seagirt city of white houses and elegant parks and gardens is no longer the capital of a vast overseas empire, it remains a busy commercial and tourist centre. Lisbon has exchanged the sounds of the past—the cries of Galician water-carriers and of bakers bearing huge baskets of bread, the whistles of knife-grinders, the bagpipes of peasants from the north—for the honking of congested motor traffic and the clang of trams.

Some traditions remain, however. One can still see in

the streets *varinas* (fish sellers), dressed in long, black skirts, carrying their wares in baskets on their heads; and one can still hear sung in the little cafés of the medieval Alfama quarter the sad, romantic music called fado. The port maintains an intimacy with its city that was common in the days before steam. Amid the freighters, warships, liners, and ferryboats, a picturesque note is struck by the *fragatas*, of Phoenician origin: these crescent-shaped boats with their striking black hulls and pink sails still perform most of the harbour's lighterage. Vessels tie up at quays open to the everyday life of the town, where the clang of the trolley cars blends with the sound of ships' bells. At dawn, fishing smacks deposit their catch at the town's front doorstep for noisy auction to Lisbon dealers, while the *varinas* wait to fill the baskets they peddle through the streets. Farther within, the fish market gives way to the equally colourful and clamorous fruit and vegetable market. Despite modernization, Lisbon, in many ways, retains the air of a 19th-century city.

THE LANDSCAPE

The city site. Lisbon is built in a succession of terraces up the slopes of a range of low, rolling hills, which rise from the banks of the Tagus River and the Sea of Straw northwest toward the Sintra Mountains, whose covering of lush Mediterranean and north European flora provides an attractive retreat for the city's population. Sections of the city vary considerably in height, especially in the older areas along the water's edge, offering splendid views of the river and the low cliffs that line the river's southern shore.

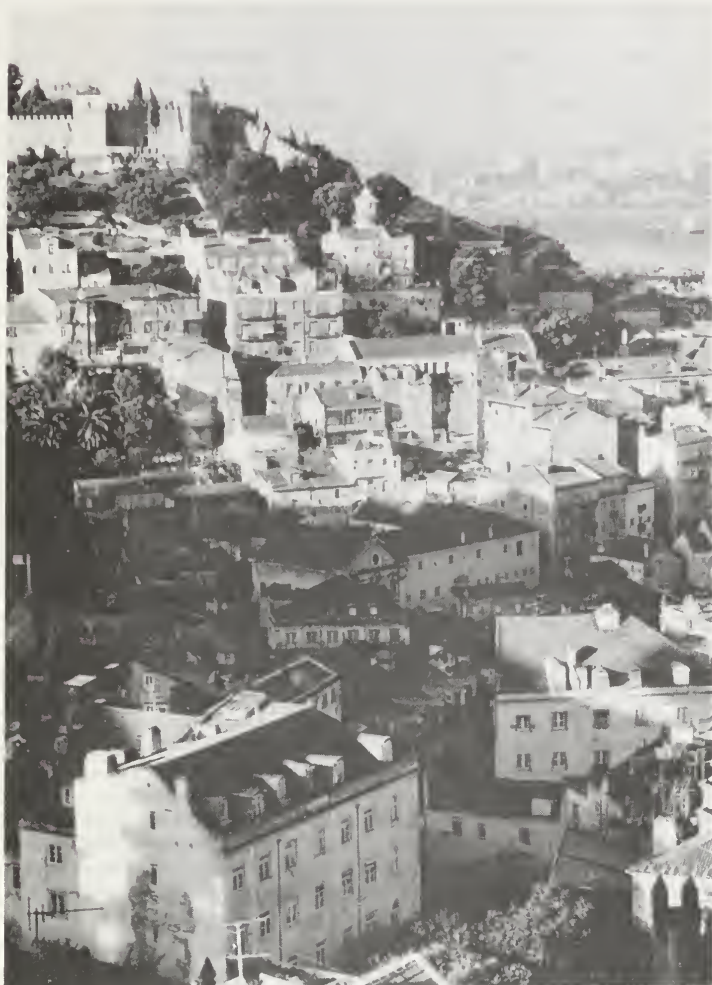
Several geologic faults cross Lisbon and the surrounding region. A major earthquake devastated the city in 1755, but seismic activity in the 20th century has been limited to slight tremors.

Climate. Lisbon has a mild and equable climate, with a mean annual temperature of 63° F (17° C). The proximity of the Atlantic and the frequency of sea fogs keep the atmosphere humid, and summers can be somewhat oppressive, although the city has been esteemed as a winter health resort since the 18th century. Average annual rainfall is 26.6 inches (666 millimetres).

The city layout. It is traditional for Lisbon's poets to refer to the entwining Tagus as the city's lover. The river is indeed an ever-present part of the city's decor, and the official entrance to Lisbon is a broad marble stair mounting from the water to the vast, arcaded Commerce Square (Praça do Comércio). The three landward sides of the square are surrounded by uniform 18th-century buildings, sea-green and white. This formal, Baroque-flavoured composition is pierced by a monumental archway, built a century later, marking the entry into the central city. In the middle of the square, surrounded each day by a regiment of parked automobiles, stands a bronze statue of King Joseph (José) I, on horseback, from which derives the nickname given the area by English sailors, Black Horse Square.

The square lies at the south end of Lisbon's central district, the Cidade Baixa ("Lower City"). The Baixa area was completely rebuilt after the earthquake in 1755 under the supervision of Joseph I's prime minister, Sebastião de Carvalho, later the marquês de Pombal. The streets are laid out in a grid pattern broken by spacious squares. A series of parallel streets, each named for its original intended occupants (e.g., Rua Áurea ["Golden Street"] for the goldsmiths), run north from Commerce Square to Dom Pedro IV Square, commonly known as Rossio Square. The Rossio is a traditional centre of activity and the starting point of the city's main promenade, the wide, gently sloping Avenida da Liberdade. This tree-lined boulevard leads north from the city centre to more modern sections of the town.

The role
of the
river



Central Lisbon and the Tagus River, with the Castle of St. George on the hilltop.

Robert Frerck/Odyssey Productions

To the east of the Baixa lies the Alfama, the oldest part of the city, where narrow, winding streets crowd down to the river between a jumble of houses. In this area, from the hill where Lisbon was first founded, the Castle of St. George (São Jorge) watches over the city. The castle, like most of the buildings in the Alfama, is Moorish in origin; it was named for England's patron saint, in honour of an alliance made in 1386. Just below it, the austere, white Church of St. Vincent-Outside-the-Walls guards the remains of the saint, which were—according to legend—miraculously brought to the city in a ship guided by two ravens. To commemorate the event, the birds are depicted on the Lisbon coat of arms.

A number of neighbourhoods extend west of Cidade Baixa toward the suburb of Belém. Each possesses its own distinctive character, reflecting the epoch in which it was built. The Bairro Alto ("Upper District"), for example, dates primarily from the 17th century. It is characterized by its straight, narrow, steeply inclined streets. Some of its streets, especially those leading down to the Baixa, are so steep that they give way to stairs, cable cars, and even an elevator (an iron structure designed by the French engineer Gustav Eiffel).

Despite new construction, the general outlines of the city remain as they were. It is still a city of balconies and vistas. On 17 of its prominences (many Lisboans profess to see only seven traditional hills, as in Rome) the *miradouros*, which are garden balconies maintained by the municipality, are still frequented by citizens of all ages.

Housing. New housing projects, hotels, and offices have begun to change the city. The pastel-tinted and somewhat sleepy Lisbon that offered a neutral, 19th-century haven to 200,000 war refugees in the 1940s has disappeared in the din and dust of new construction. Lisbon has emerged as a bustling modern metropolis.

Pombal's Baixa remains rigorously protected from change, but the four-story buildings of the Avenida da Liberdade and its ancillary streets have been almost totally replaced by 10-story buildings in a bland modern style. The new construction has also gained the hills, even the Alfama, whose established residents continue to hang laundry across the narrow streets and to grill sardines on doorstep braziers.

The municipality has built new neighbourhoods in the northern and northwestern sectors of town. Other developments have pushed westward toward Belém as a replacement for decayed neighbourhoods. These modern structures, some of them 14 stories high, are designed by Lisbon architects, who produce handsome, colourful, contemporary buildings. Most of the new lodgings are reserved for the poorest families (a certain number pay nothing at all) and for people of moderate means. The Chelas District project, implanted on heathland previously considered too difficult to build upon, houses about 10 percent of the Lisbon population. Despite these projects, adequate housing remains a problem, and a number of shantytowns have developed on the periphery of the city. Many affluent families have moved out of town, and more and more villas have appeared in the countryside of the "Portuguese Riviera," which lies between Lisbon and the town of Estoril, 16 miles to the west.

New housing

THE PEOPLE

Although the district of Lisbon occupies only about 3 percent of Portugal's total area, it contains more than 20 percent of the country's population. Similarly, the city proper occupies approximately 3 percent of Lisbon district, but roughly 40 percent of the district's residents live in the city. The area has long been a magnet for rural immigration. The population is predominantly Portuguese, but there are some foreign residents, mostly diplomats and merchants in the import-export business. Lisbon, judged by net median family income, is the capital of western Europe's poorest people.

For centuries the Lisboans have discussed the symptoms of an affliction said to be endemic in this strip of the Iberian Peninsula: *saudade*, generally translated as "melancholy," a variety of a state of anxiety tempered by fatalism. *Saudade* is said to be reflected in the fado, which is sung in its original form only in Lisbon, and in Lisbon only in the two hillside precincts of Alfama and Bairro Alto. The word fado means "fate," usually an unkind fate described in the songs that are throbbingly sung to a penetrating but melodic music.

Saudade

Portugal is an essentially Roman Catholic country, and Lisbon is distinguished as one of the three places in the world whose chief Roman Catholic clergyman bears the title of patriarch. The Lisboans are typically less devout than northern Portuguese, however, using the church mainly for family occasions: christenings, weddings, and funerals. Religious processions are generally subdued affairs, without the colour and the drama found in Spain. The June feasts of the popular saints (St. Anthony, St. John, and St. Peter) are exceptions. The Lisboans celebrate by donning imaginative costumes, jumping over bonfires, and dancing in the streets until dawn. Indeed, these lively events, held in the small quarters near the Castle of St. George, retain all the pagan elements of a midsummer festival.

THE ECONOMY

Industry. To the long-established local industries of soapmaking, munitions, and steel manufacture have been added glassmaking, electronics, margarine manufacture, and diamond cutting. The petroleum refinery, largely state supported, has also been expanded.

The greatest development in recent years, however, has come just outside the city limits on the south bank of the Tagus; the district has become Portugal's most important manufacturing centre. The industrial belt has continued to grow from the river down to Setúbal, which is 25 miles south. One of the world's largest cement plants is found on the far bank, along with grain elevators, a steelmaking complex, a cork factory, and a plastics plant.

The port of Lisbon, with 19 miles of docks, now has

The role of tourism

special facilities for the handling of container-ship cargoes and for car ferries, and it continues to undergo expansion. Cotton, grain, and coal are important imports.

Commerce and finance. Tourism and commerce have played a major part in Lisbon's modernization, and revenues from tourism have helped offset usually negative national trade balances. With the increase in population and the spread of hotels, offices, and apartment blocks throughout the city, banks have proliferated. Shops are concentrated around the Rua Garrett and the Baixa, though good shopping centres have grown up in the new residential districts toward the airport. The city's major market is located in the riverside square near the Cais do Sodré railway station.

The commercial docks, situated in Alcântara, west of Cidade Baixa, have refrigeration plants located close by to handle the catches of sardines and tunny on which the trawler fleet's livelihood has long been based.

Transportation. Lisbon is connected by rail and road to the interior of Portugal and to the rest of Europe. The airport, at Portela de Sacavém, some six miles beyond the city, has flights to Europe, the Americas, Africa, and the Middle East. The 25th of April Bridge has been the main roadway into the city since it was built in the mid-1960s. The first bridge in Lisbon's long history to span the Tagus, it is more than 7,470 feet (2,277 metres) from anchorage to anchorage, with a central span of about 3,323 feet suspended 230 feet above mean water level. There is space under the roadway to carry two railroad tracks.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. As the capital of Portugal, Lisbon and its surrounding suburbs house all of the principal institutions of the republic. Many government offices occupy the 18th-century buildings that surround Commerce Square. The nation's parliamentary body meets at the Palace of the National Assembly (Palácio da Assembléia Nacional), which also houses the National Archives. Located on the west side of the Bairro Alto, the 17th-century buildings were originally occupied by the convent of São Bento da Saúde. Farther west, toward Belém, the Necessidades Palace (Palácio das Necessidades) houses the Ministry of Foreign Affairs. Built in the mid-18th century on the site of a chapel (for which it is named), it was, until 1910, a royal palace. Another royal palace, built in Belém in 1700, is now the official residence of the president of Portugal.

Like Portugal's 21 other administrative districts, the district of Lisbon is headed by a civil governor who is appointed by and responsible to the central government through the Minister of the Interior. The district is divided into municipalities (*concelhos*), which are further subdivided into wards (*bairros*) and parishes (*freguesias*). Representatives to parish assemblies are chosen through local elections. They then elect an executive body, the parish committee. Lisbon's Municipal Assembly consists of representatives chosen by their parish committees and members directly elected by the local citizens. It serves as the legislative branch of local government and elects the executive branch, the Municipal Chamber, which is headed by a president appointed by the district governor. Services under the jurisdiction of the municipality include the city waterworks, road maintenance, and sanitation. The money for these services is provided by grants from the central government and through local taxes.

Health. As a major urban centre, Lisbon has a higher percentage of doctors and other health professionals than the rest of Portugal. The city's hospitals include state, private, and military establishments. There are no hospitals run by the municipality. State hospitals are handicapped by the scarcity of public funding. As a result they are insufficiently staffed and, in some cases, inadequately equipped to cope with the demands put on their services. English, French, and Hebrew hospitals cater largely to their small, respective communities.

Education. Officially, education in Portugal is free and compulsory for children ages six through 14, but the number of public schools in Lisbon is insufficient for the city's population. Private schools fill the gap to a certain extent. The University of Lisbon was founded in 1290 and re-

mained Portugal's only university until the 16th century. It was moved back and forth between Lisbon and Coimbra until 1537, when it was renamed and permanently located in Coimbra. Thus, Lisbon was left without a university until 1911, when a second University of Lisbon was established. The Technical University of Lisbon was founded in 1931, and three more universities were opened in the city during the 1970s. Despite this effort to expand educational opportunities, the number of university applicants far exceeds the number of available places.

CULTURAL LIFE

Lisbon's Calouste Gulbenkian Foundation and Museum, a cultural centre named for its benefactor, an Armenian oil-lease negotiator, presents music and ballet as well as exhibits of the other fine arts. It also houses the broad-ranging personal collection that Calouste Gulbenkian, who lived in Lisbon from 1942 until his death in 1955, willed to the Portuguese nation.

The city has more than a dozen other museums, including those of modern, antique, sacred, decorative, and folk arts. In the latter are found many beautiful native artifacts. Two specialized, rather unusual museums are the Azulejo Museum and the Coach Museum. The former, located in the convent of Madre de Deus, boasts a large and varied collection of the painted tiles (*azulejos*) for which the Iberian Peninsula is famous. The Coach Museum occupies a wing of the Portuguese president's official residence and contains an impressive display of carved and gilded coaches.

Lisbon's municipal orchestra was founded in 1971. The National Conservatory offers advanced instruction in both music and drama. The city has two principal theatres, the St. Charles (São Carlos) and the National Theatre of Dona Maria II. The St. Charles, which was constructed in the late 18th century, has a beautiful elliptical interior, and the National, which was built about 1845, displays a facade of six giant columns saved from the convent church of St. Francisco, which was destroyed by earthquake. The interior, gutted by fire in 1966, has been restored.

Neither of these edifices is as theatrical as the interiors of some of the churches built or restored after the 1755 earthquake. In gold, marble, carved wood, and rare tiles, these interiors are decorated in Baroque, Rococo, or rocaille. One outstanding example is the 16th-century church of São Roque, whose unpretentious exterior belies its opulent collection of *azulejos*, paintings, and mosaics inlaid with semiprecious stones.

The old, red brick bullring with its Moorish arches and cupolas still finds an audience for its spectacles. In Portuguese bullfights the bull's first opponent is either a *cavaleiro* on horseback or a *toureiro* on foot. The performance of either man is judged by his dexterity, courage, and proximity to the bull. Next, a squad of acrobatic men, the *forcados*, wrestle the bull, eventually immobilizing it with their bare hands. The bull is not slain.

Lisbon has several other sports and recreational areas. Residents can travel several miles north or west to one of the three major football (soccer) stadiums. Many of the housing developments are planted with trees and grass, their small parks adding to Lisbon's collection of more than 40 public gardens. The largest public park, Monsanto, covers about eight square miles and has numerous recreational facilities. Its rolling hills were planted in the 1920s to provide a windbreak for the town and are now thickly forested. There are also two botanical gardens and a zoological garden within the city.

History

THE EARLY PERIOD

Prehistoric to Moorish times. The valley in which the heart of Lisbon now lies was, in prehistoric times, the bed of a forked branch of the Tagus. (The subway now forks at the same spot.) No evidence has been uncovered to show who were the first residents on the hills surrounding the valley. Although it seems likely that the city was founded c. 1200 BC as a trading station by the far-ranging Phoenicians, there is no unassailable proof of the story.

Theatres and theatrical churches

Municipal services

Origins of the city

The city's ancient name, *Olisipo*, may be derived from the Phoenician *alis ubbo* ("delightful little port") or from the legend that the city's founder was Ulysses.

Whatever the city's origins, it is known that the area was under Roman domination from 205 BC to c. AD 409 and that Julius Caesar raised the settlement to the dignity of a *municipium* and named it *Felicitas Julia*. A few inscribed stones remain as evidence of the Roman presence. The Romans lost the city to the migratory peoples known as the Alani, who were driven out by the Suebi, who in turn were conquered by the Visigoths. The base plan of the original fortifications is thought to be Visigothic and, if so, is the sole vestige of their reign.

The Muslims of North Africa took Lisbon when they overran the Iberian Peninsula in the 8th century; they stayed for 433 years, despite incursions by the Normans in 844 and by Alfonso VI of Castile and León in 1093. Under the Moors the city was known under variations of "Lisbon": *Luzbona*, *Lixbuna*, *Ulixbone*, and *Olissibona*. Some authorities contend that the Muslims took this name from the conquered Roman castle, but Lisbon historians suggest that it derives from *água boa* ("good water").

The Portuguese conquest. Behind their walls, the Moors were able to hold out for months when the city was assailed by crusader forces—English, Flemish, Norman, and Portuguese under Afonso Henriques, the Portuguese king. The city finally fell in 1147 and successfully resisted Moorish attempts to win it back. The Moorish alcazar was transformed into a Portuguese royal palace, and, according to legend, the Lisbon Cathedral (*Sé Patriarcal*) was converted from a mosque (with subsequent restorations in the styles of many periods after fires and earthquakes). There is no evidence, however, of a building on the site of the cathedral before the time of Afonso Henriques.

Although 1,400 years of occupation and invasion have left almost no trace among the stones of the capital, the presence of the outlander is still visible in the faces of the inhabitants, which range in cast from the Scandinavian to the Mauritanian.

After winning Lisbon, King Afonso established his court 105 miles to the north-northeast, atop a cliff at Coimbra. Lisbon did not become the national capital until more than a century later, in 1256. Within its Moorish walls, of which some traces still remain, medieval Lisbon measured 1,443 feet at its widest and 1,984 feet at its longest, descending the hill below the castle. Even before the Portuguese conquest, some houses had already been built outside the walls toward the river. The site of this first Lisbon is occupied by the lively Alfama quarter, which has kept the labyrinthine medieval street plan.

King Dinis I (1279–1325) decreed that Portuguese, the dialect of the Porto region, was to be the national language. He founded the university in Lisbon in 1290, and during his reign, other hilltops around the central valley were crowned with convents and churches.

In 1372–73 Lisbon was besieged and burned by the Castilians, who forced King Ferdinand I, an unsuccessful contender for the Castilian throne, to repudiate his alliance with England; thereafter the King swiftly erected new defenses. His wall—more than three miles long, with 77 towers and 38 gates and enclosing more than 247 acres—withstood the renewed Castilian attack of 1384, which followed Ferdinand's death.

The Age of Discovery. When the Portuguese Age of Discovery (1415–1578) began, a census of Lisbon showed 65,000 inhabitants occupying 23 parishes. A considerable number of these residents became rich, and the city was endowed with larger and more luxurious buildings. African slaves became a familiar Lisbon sight, the trade in slaves being one in which Portugal played a major role. After the great explorer Vasco da Gama led a Portuguese fleet to India in 1498, the Venetian monopoly on Oriental trade was broken; and colonies of German, Flemish, Dutch, English, and French traders established themselves in Lisbon. Greeks, Lombards, and Genoese who had lost their trading enclaves in Constantinople when that city fell to the Turks in 1453 also came to Lisbon.

King Manuel I (1495–1521) dominated this epoch, and under his rule Portugal developed its sole contribution to

European architecture, an extreme style of late Gothic decoration that celebrated the voyages of discovery. Manuel, and God. The prime examples of Manueline style at Lisbon, the Tower of Belém and the Jerónimos Monastery, about four miles downstream from the city centre, are far less exuberant than those at the rival Portuguese cities of Batalha and Tomar. The tower and the monastery are nevertheless the most important architectural monuments in the Lisbon area. The five-story Tower of Belém, located on the riverbank, was built in 1515 as a fort in the middle of the Tagus, which subsequently altered course. Girt by a cable carved in the stone, it has a stern Gothic interior but exhibits its North African touches on its turrets and crenellations and presents rounded Renaissance arches for the windows. The monastery with its church and cloisters was begun in 1502 by Boytac (Boitaca), an architect of French origin, and was not finished until the end of the century. Four other architects worked on the project, their styles passing from the Gothic through the Renaissance to the Baroque. Smoothed by time, the ensemble is harmonious and proudly Portuguese.

Manuel I promoted the urbanization of the central valley between Lisbon's hills, creating a city square, the Rossio, which at once became a popular meeting place. By the Tagus he constructed a new palace, the Paços da Ribeira, with a large square laid out along its eastern flank. The area between the Rossio and the Palace Terrace (*Terreiro do Paço*) was soon crisscrossed with streets, along which rose the new shops, churches, and hospitals of what had become a phenomenally prosperous city.

The prosperity was chimerical, however. John (João) the Pious, who had succeeded Manuel, permanently transferred (1537) the university to the royal palace at Coimbra, far from the capital's excesses. He also invited the Jesuits and the Inquisition to come to Portugal to counter the ungodly materialism of Lisbon. The Inquisition office, located in the Rossio, was particularly ferocious in its persecution of the Jews, who were the bankers, financiers, and moneylenders of the time. Many wealthy Jews had their property and goods confiscated; some emigrated to Holland and other countries, taking their money and financial expertise with them. As a result, Lisbon's connections with foreign markets were disrupted and the country's economy suffered severe financial constraints.

Lisbon was visited with plagues and earthquakes during this time, but they proved easier to meet than the cost of 50 years of glory. Literally half the nation's population had vanished in pursuit of wealth in the new colonies. With farms deserted, food was imported from other European countries at crippling prices, and with so many skilled men absent, wages rose sharply, as did the cost of building and manufacturing materials. The colonial treasures, which had made Lisbon such a sybaritic queen of the seas, in the end cost more than they could fetch.

In 1578 King Sebastian of Portugal was killed in a disastrous invasion of Morocco; two years later, the Spanish pushed into Portugal, and Philip II of Spain became king of both countries. In 1588 it was from Lisbon that the Invincible Armada sailed against England, Portugal's oldest ally. In the half century that followed, Lisbon lived relatively well as a port for the riches of the Spanish Main. In 1640 a conspiracy of Lisbon nobles struck for freedom and drove out the Spaniards, restoring Portugal's independence. The square just north of the Rossio, Restoration Square (*Praça dos Restauradores*), is named for them.

With the Cromwellian treaty of 1654, following British military assistance to the Portuguese in the war with Spain, the British merchants trading and living in Lisbon set up a corporation, which became known as the British Factory. The Factory negotiated with the Portuguese government for trade concessions and other privileges, appealing to the British government to put pressure on the Portuguese authorities when necessary. Britain's economic and political influence on Portugal was strong, and the Factory remained in existence until 1810.

EVOLUTION OF THE MODERN CITY

Disaster and reconstruction. In the first half of the 18th century the profits from the plantations and the gold

Decline in the 16th century



The Jerónimos Monastery in Lisbon.
Colour Library International

The great earthquake of 1755

and diamond deposits of Brazil brought a new flurry of optimism and excitement to Lisbon. This ended on the morning of Nov. 1, 1755. The churches were crowded to honour the dead on All Saints' Day when the city was devastated by one of the greatest earthquakes ever recorded. There were two shocks, 40 minutes apart; and the waters of the Tagus, lifted from their bed, roared through the city, followed by fire. It is believed that 30,000 lives were lost, and more than 9,000 buildings were destroyed.

Physically, Lisbon recovered with a celerity astonishing for the time, but the shock left its mark upon the thinking of generations to come. The reconstruction—a good deal of foreign aid was forthcoming—was achieved by Joseph I's prime minister, Sebastião José de Carvalho, the virtual ruler of the realm, in charge of five architects and soon had a plan for remaking the totally devastated centre of the town, the Baixa. The riverside palace had been destroyed, and its terrace was expanded to create the new Commerce Square. Northward from there, a grid of 48 streets led inland to the Rossio and a neighbouring new square, Figueira. The two-story, uniform buildings were topped by two tiers of dormers projecting from tiled roofs. The corners of the eaves, in the Lisbon tradition, turned up, in faint echo of a pagoda. The building style, evolved for fast, cheap construction, was Baroque but virtually stripped of decoration. After the minister was rewarded with the title of *marquês de Pombal*, the style became known as the *estilo pombalino*.

The Sé and most of the churches were repaired or rebuilt, but the 14th-century Carmel (Carmo) Church was left as it was. Looming from its hilltops over the Baixa, the roofless Gothic shell now serves as an archaeological museum, while its cloister serves as the barracks for the National Republican Guard, a paramilitary security force. The Palace of the Inquisition, utterly flattened, was not rebuilt when Pombal enlarged and realigned the Rossio, and on its site, 90 years later, the National Theatre of Dona Maria II was erected. Pombal banished the Jesuit order and transformed their establishment into St. Joseph's Hospital to replace the destroyed All Saints Hospital. The medical school scrambled for room at St. Joseph's until it acquired a new building of its own late in the 19th century. The Jesuit novice house was converted to serve as the Nobles' School. Later governments expelled more religious orders, whose buildings became barracks, hospitals, royal academies, and government offices.

19th-century expansion. During the Peninsular War of the early 1800s, Lisbon alternated between French and British control, and after Napoleon's defeat it was embroiled in civil war until 1834. That conflict was followed by 10 years of revolutionary outbursts. Nineteenth-century Lisbon nevertheless continued to expand and, by 1885, embraced 20,378 acres, while the population had doubled

in 100 years to reach 300,000. Public buildings, such as the new city hall and the Ajuda Royal Palace, had been built, and the harbour had been modernized and quays constructed on land reclaimed from the river. The railway had appeared, and a system of horse cars served the lower town.

The greatest change in the city, and the one most important for its future growth, was the opening of a new main street in 1880—the Avenida da Liberdade. The municipality bordered the central six-lane carriageway with wide blue mosaic sidewalks graced with palms and shade trees, fountains, and ornamental waters stocked with goldfish and swans. So the street remains today, with the addition of outdoor cafés beneath the trees.

In conjunction with the new thoroughfare, a series of new streets, the "Avenidas Novas," expanded the city northward, and new neighbourhoods sprang up as well on the borders of the Avenida da Liberdade. In 1901 the electric streetcar made its appearance, enabling more people to live farther away from their employment in the Baixa. Three cable cars shuttled up and down the adjacent hills, and Eiffel designed the giant elevator that hisses grandly between the town's upper and lower levels.

New water supplies, augmenting those of the 1748 aqueduct of *Águas Livres*, were introduced from Alviela. Consequently, water was piped directly into houses, eliminating the ancient calling of *galego*, or water porter.

The 20th century. In 1908 Portugal's king and crown prince were assassinated on the northwest corner of Commerce Square. Two years later the new king, Manuel II, abdicated. A republic was declared, and a period of national instability ensued. When António Salazar took control of the near-bankrupt nation in 1932, he erected a corporate state of which he alone determined the policies until his retirement in 1968. There was considerable development in Lisbon throughout this time, and during World Wars I and II, in which Portugal was neutral, the city was able to offer refuge to some 200,000 foreigners.

In the 1960s national policy began to change, allowing economic expansion. The 30-year-old austerity program of stability and self-sufficiency (at an admittedly low level of investment and consumption) was somewhat softened, and international tourists and international corporations began to be accommodated. In 1966, well ahead of schedule (and with \$77,000,000 in aid from the United States), the Salazar (now the 25th of April) Bridge was completed.

On April 25, 1974, the government of Marcelo Caetano, Salazar's successor, was overthrown by a military coup. By the early 1980s, however, political instability and economic difficulties remained serious problems and hindered the nation's, and the city's, efforts to bring about social and economic reforms. (B.E./L. de S.R.)

In 1988 a fire destroyed much of the city's historic Chiado district, which was rebuilt into a shopping and cultural centre during the 1990s. Lisbon was designated a European City of Culture in 1994. The city's landscape was altered in 1998, when Lisbon hosted the World's Fair (Expo '98), which sparked the biggest renewal project since the 1755 earthquake and tidal wave. New facilities included the cable-stayed Vasco Da Gama Bridge, the Oceanário de Lisboa (Europe's largest oceanarium), and marinas, hotels, and commercial complexes along the Tagus. (Ed.)

BIBLIOGRAPHY

General: Fodor's *Lisbon*, 1986 (1986), provides general descriptive information. DAVID WRIGHT and PATRICK SWIFT, *Lisbon: A Portrait and a Guide* (1971), gives thorough coverage of all quarters of the city, its history, monuments, cultural institutions, and contemporary life and also includes excursions outside Lisbon. Also see CAROL WRIGHT, *Lisbon* (1971), on the city and the life of its inhabitants at various times of day, as well as information for the visitor to Lisbon and its environs; and VIVIAN ROWE, *The Road to Lisbon* (1962), mainly concerned with the journey from France to Lisbon but including some material on the attractions of the city itself.

History: JÚLIO DE CASTILHO, *Lisboa Antiga: O Bairro Alto de Lisboa*, 2nd ed., 5 vol. (1902–04), and *Lisboa Antiga Bairros Orientais*, 2nd ed., 12 vol. (1934–38); THOMAS D. KENDRICK, *The Lisbon Earthquake* (1956), a study of the 1755 earthquake and its impact on European philosophers and theologians.

(B.E./L. de S.R.)

The opening of Avenida da Liberdade

An era of expansion

The Art of Literature

The current meaning generally attached to the term literature—a body of writing by a people or by peoples using the same language—is a relatively modern one. The term itself, derived from the Latin word *littera* (“letter of the alphabet”; *litterae*, “letters”), is ancient enough; but in ancient times literature tended to be considered separately in terms of kinds of writing, or genres as they came to be called in the 18th century when the term literature took on its modern meaning. Thus Aristotle’s *Poetics*, though it is concerned with and gives examples from Greek epic and dithyrambic poetry and comedy, has as its central concern tragedy.

The sections that make up this article are arranged roughly in a chronological order according to the mediums of verse and prose with attention, again mainly chronological, to various genres. No such classification is totally satisfactory. Furthermore, although the derivation of the word literature implies writing, there is much oral literature, a general treatment of which may be found in the article FOLK ARTS: *Folk literature*.

There is a further complication: literature as a whole and in its parts means various things to various writers, critics, and historians. At one extreme, it may be held that any-

thing written is literature. Though this position is seldom held, that at the other extreme—literature is only the *Iliad*, the *Odyssey*, and *Hamlet*—is slightly more popularly held. Between these extremes, attitudes vary widely. For some critics, a hierarchy exists: tragedy is superior to comedy; the short story is inferior to the novel. For other critics, qualitative criteria apply: poetry is verse that succeeds; the limerick and nonsense verse are failed poetry. Critics also differ on the purpose or ends of literature. Many ancient critics—and some modern ones—hold that the true ends of literature are to instruct and delight. Others—a majority of the modern ones, probably—hold that pleasure is the sole end. All of these divergences and other similar ones appear in the treatments that follow.

For historical treatment of various literatures, see the article LITERATURE, THE HISTORY OF WESTERN, and sections in the articles AFRICAN ARTS, CENTRAL ASIAN ARTS, OCEANIC ARTS, and SOUTHEAST ASIAN ARTS. Some literatures are treated separately by language, by nation, or by special subject (e.g., CELTIC LITERATURE, LATIN LITERATURE, FRENCH LITERATURE, JAPANESE LITERATURE, BIBLICAL LITERATURE).

The article is divided into the following sections:

-
- The scope of literature 78
 - Literary composition
 - Content of literature
 - Literature and its audience
 - Literature and its environment
 - Literature as a collection of genres
 - Writings on literature
 - Poetry 85
 - The nature of poetry 85
 - Attempts to define poetry
 - Poetry and prose
 - Form in poetry
 - Poetry as a mode of thought: the Protean encounter
 - Prosody 89
 - Elements of prosody
 - Prosodic style
 - Theories of prosody
 - Narrative fiction 96
 - Epic 96
 - General characteristics
 - Bases
 - Early patterns of development
 - Later variations
 - Fable, parable, and allegory 101
 - Nature and objectives
 - Historical development in Western culture
 - Allegorical literature in the East
 - Ballad 107
 - Elements
 - Composition
 - Types of balladry
 - Subject matter
 - Chronology
 - Romance 110
 - The component elements
 - Medieval verse romances
 - Medieval prose romances
 - Later developments
 - Saga 114
 - Nonfictional saga literature
 - Legendary and historical fiction
 - Novel 116
 - Elements
 - Uses
 - Style
 - Types of novel
 - The novel in English
 - Europe
 - Asia, Africa, Latin America
 - Social and economic aspects
 - Evaluation and study
 - The future of the novel
 - Short story 138
 - Analysis of the genre
 - History
 - The 20th century
 - Drama 143
 - Dramatic literature 143
 - General characteristics
 - Drama as an expression of a culture
 - Influences on the dramatist
 - The range of dramatic forms and styles
 - Comedy 151
 - Origins and definitions
 - Theories
 - Kinds of comedy in diverse historical periods
 - The comic in other arts
 - Tragedy 160
 - Development
 - Tragedy and modern drama
 - Theory of tragedy
 - Other genres 173
 - Satire 173
 - The nature of satire
 - Satirical media
 - The satirist, the law, and society
 - Nonfictional prose 176
 - Nature
 - Elements
 - Approaches
 - The essay
 - History
 - Doctrinal, philosophical, and religious prose
 - Political, polemical, and scientific prose
 - Other forms
 - Biographical literature 186
 - Aspects
 - Kinds
 - Historical development
 - Biographical literature today
 - Literary criticism 194
 - Functions
 - Historical development
 - The 20th century
 - Children’s literature 198
 - Definition of terms
 - The case for a children’s literature
 - Some general features and forces
 - The development of children’s literature
 - Historical sketches of the major literatures
 - Bibliography 211
-

THE SCOPE OF LITERATURE

Literature is a form of human expression. But not everything expressed in words—even when organized and written down—is counted as literature. Those writings that are primarily informative—technical, scholarly, journalistic—would be excluded from the rank of literature by most, though not all, critics. Certain forms of writing, however, are universally regarded as belonging to literature as an art. Individual attempts within these forms are said to succeed if they possess something called artistic merit and to fail if they do not. The nature of artistic merit is less easy to define than to recognize. The writer need not even pursue it to attain it. On the contrary, a scientific exposition might be of great literary value and a pedestrian poem of none at all.

The purest (or, at least, the most intense) literary form is the lyric poem, and after it comes elegiac, epic, dramatic, narrative, and expository verse. Most theories of literary criticism base themselves on an analysis of poetry, because the aesthetic problems of literature are there presented in their simplest and purest form. Poetry that fails as literature is not called poetry at all but verse. Many novels—certainly all the world's great novels—are literature, but there are thousands that are not so considered. Most great dramas are considered literature (although the Chinese, possessors of one of the world's greatest dramatic traditions, consider their plays, with few exceptions, to possess no literary merit whatsoever).

The Greeks thought of history as one of the seven arts, inspired by a goddess, the muse Clio. All of the world's classic surveys of history can stand as noble examples of the art of literature, but most historical works and studies today are not written primarily with literary excellence in mind, though they may possess it, as it were, by accident.

The essay was once written deliberately as a piece of literature: its subject matter was of comparatively minor importance. Today most essays are written as expository, informative journalism, although there are still essayists in the great tradition who think of themselves as artists. Now, as in the past, some of the greatest essayists are critics of literature, drama, and the arts.

Some personal documents (autobiographies, diaries, memoirs, and letters) rank among the world's greatest literature. Some examples of this biographical literature were written with posterity in mind, others with no thought of their being read by anyone but the writer. Some are in a highly polished literary style; others, couched in a privately evolved language, win their standing as literature because of their cogency, insight, depth, and scope.

Many works of philosophy are classed as literature. The *Dialogues* of Plato (4th century BC) are written with great narrative skill and in the finest prose; the *Meditations* of the 2nd-century Roman emperor Marcus Aurelius are a collection of apparently random thoughts, and the Greek in which they are written is eccentric. Yet both are classed as literature, while the speculations of other philosophers, ancient and modern, are not. Certain scientific works endure as literature long after their scientific content has become outdated. This is particularly true of books of natural history, where the element of personal observation is of special importance. An excellent example is Gilbert White's *Natural History and Antiquities of Selbourne* (1789).

Oratory, the art of persuasion, was long considered a great literary art. The oratory of the American Indian, for instance, is famous, while in classical Greece, Polymnia was the muse sacred to poetry and oratory. Rome's great orator Cicero was to have a decisive influence on the development of English prose style. Abraham Lincoln's Gettysburg Address is known to every American schoolchild. Today, however, oratory is more usually thought of as a craft than as an art. Most critics would not admit advertising copywriting, purely commercial fiction, or cinema and television scripts as accepted forms of literary expression, although others would hotly dispute their exclusion. The test in individual cases would seem to be one of endur-

ing satisfaction and, of course, truth. Indeed, it becomes more and more difficult to categorize literature, for in modern civilization words are everywhere. Man is subject to a continuous flood of communication. Most of it is fugitive, but here and there—in high-level journalism, in television, in the cinema, in commercial fiction, in westerns and detective stories, and in plain, expository prose—some writing, almost by accident, achieves an aesthetic satisfaction, a depth and relevance that entitle it to stand with other examples of the art of literature.

LITERARY COMPOSITION

Critical theories. *Western.* If the early Egyptians or Sumerians had critical theories about the writing of literature, these have not survived. From the time of classical Greece until the present day, however, Western criticism has been dominated by two opposing theories of the literary art, which might conveniently be called the expressive and constructive theories of composition.

The Greek philosopher and scholar Aristotle (384–322 BC) is the first great representative of the constructive school of thought. His *Poetics* (the surviving fragment of which is limited to an analysis of tragedy and epic poetry) has sometimes been dismissed as a recipe book for the writing of potboilers. Certainly, Aristotle is primarily interested in the theoretical construction of tragedy, much as an architect might analyze the construction of a temple, but he is not exclusively objective and matter of fact. He does, however, regard the expressive elements in literature as of secondary importance, and the terms he uses to describe them have been open to interpretation and a matter of controversy ever since.

The 1st-century Greek treatise *On the Sublime* (conventionally attributed to the 3rd-century Longinus) deals with the question left unanswered by Aristotle—what makes great literature "great"? Its standards are almost entirely expressive. Where Aristotle is analytical and states general principles, the pseudo-Longinus is more specific and gives many quotations: even so, his critical theories are confined largely to impressionistic generalities.

Thus, at the beginning of Western literary criticism, the controversy already exists. Is the artist or writer a technician, like a cook or an engineer, who designs and constructs a sort of machine that will elicit an aesthetic response from his audience? Or is he a virtuoso who above all else expresses himself and, because he gives voice to the deepest realities of his own personality, generates a response from his readers because they admit some profound identification with him? This antithesis endures throughout western European history—Scholasticism versus Humanism, Classicism versus Romanticism, Cubism versus Expressionism—and survives to this day in the common judgment of our contemporary artists and writers. It is surprising how few critics have declared that the antithesis is unreal, that a work of literary or plastic art is at once constructive and expressive, and that it must in fact be both.

Eastern. Critical theories of literature in the Orient, however, have been more varied. There is an immense amount of highly technical, critical literature in India. Some works are recipe books, vast collections of tropes and stylistic devices; others are philosophical and general. In the best period of Indian literature, the cultural climax of Sanskrit (c. 320–490), it is assumed by writers that expressive and constructive factors are twin aspects of one reality. The same could be said of the Chinese, whose literary manuals and books on prosody and rhetoric are, as with the West, relegated to the class of technical handbooks, while their literary criticism is concerned rather with subjective, expressive factors—and so aligns itself with the pseudo-Longinus' "sublime." In Japan, technical, stylistic elements are certainly important (Japanese discrimination in these matters is perhaps the most refined in the world), but both writer and reader above all seek qualities of subtlety and poignancy and look for intimations

Distinction between poetry and verse

Role of oratory in literature

Constructive and expressive theories of criticism

of profundity often so evanescent as to escape entirely the uninitiated reader.

Broad and narrow conceptions of poetry. Far Eastern literary tradition has raised the question of the broad and narrow definitions of poetry (a question familiar in the West from Edgar Allan Poe's advocacy of the short poem in his "Poetic Principle" [1850]). There are no long epic poems in Chinese, no verse novels of the sort written in England by Robert Browning or Alfred Lord Tennyson in the 19th century. In Chinese drama, apart from a very few of the songs, the verse as such is considered doggerel. The versified treatises on astronomy, agriculture, or fishing, of the sort written in Greek and Roman times and during the 18th century in the West, are almost unknown in the Far East. Chinese poetry is almost exclusively lyric, meditative, and elegiac, and rarely does any poem exceed 100 lines—most are little longer than Western sonnets; many are only quatrains. In Japan this tendency to limit length was carried even further. The ballad survives in folk poetry, as it did in China, but the "long poem" of very moderate length disappeared early from literature. For the Japanese, the *tanka* is a "long poem": in its common form it has 31 syllables; the *sedōka* has 38; the *dodoitsu*, imitating folk song, has 26. From the 17th century and onward, the most popular poetic form was the haiku, which has only 17 syllables.

This development is relevant to the West because it spotlights the ever-increasing emphasis which has been laid on intensity of communication, a characteristic of Western poetry (and of literature generally) as it has evolved since the late 19th century. In the Far East all cultivated people were supposed to be able to write suitable occasional poetry, and so those qualities that distinguished a poem from the mass consequently came to be valued above all others. Similarly, as modern readers in the West struggle with a "communication avalanche" of words, they seek in literature those forms, ideas, values, vicarious experiences, and styles that transcend the verbiage to be had on every hand.

Literary language. In some literatures (notably classical Chinese, Old Norse, Old Irish), the language employed is quite different from that spoken or used in ordinary writing. This marks off the reading of literature as a special experience. In the Western tradition, it is only in comparatively modern times that literature has been written in the common speech of cultivated men. The Elizabethans did not talk like Shakespeare nor 18th-century people in the stately prose of Samuel Johnson or Edward Gibbon (the so-called Augustan plain style in literature became popular in the late 17th century and flourished throughout the 18th, but it was really a special form of rhetoric with antecedent models in Greek and Latin). The first person to write major works of literature in the ordinary English language of the educated man was Daniel Defoe (1660?–1731), and it is remarkable how little the language has changed since. *Robinson Crusoe* (1719) is much more contemporary in tone than the elaborate prose of 19th-century writers like Thomas De Quincey or Walter Pater. (Defoe's language is not, in fact, so very simple: simplicity is itself one form of artifice.)

Ambiguity. Other writers have sought to use language for its most subtle and complex effects and have deliberately cultivated the ambiguity inherent in the multiple or shaded meanings of words. Between the two world wars, "ambiguity" became very fashionable in English and American poetry and the ferreting out of ambiguities—from even the simplest poem—was a favourite critical sport. T.S. Eliot in his literary essays is usually considered the founder of this movement. Actually, the platform of his critical attitudes is largely moral, but his two disciples, I.A. Richards in *Principles of Literary Criticism* (1924) and William Empson in *Seven Types of Ambiguity* (1930), carried his method to extreme lengths. The basic document of the movement is C.K. Ogden and I.A. Richards' *The Meaning of Meaning* (1923), a work of enormous importance in its time. Only a generation later, however, their ideas were somewhat at a discount.

Translation. Certainly, William Blake or Thomas Campion, when they were writing their simple lyrics, were

unaware of the ambiguities and multiple meanings that future critics would find in them. Nevertheless, language is complex. Words do have overtones; they do stir up complicated reverberations in the mind that are ignored in their dictionary definitions. Great stylists, and most especially great poets, work with at least a half-conscious, or subliminal, awareness of the infinite potentialities of language. This is one reason why the essence of most poetry and great prose is so resistant to translation (quite apart from the radically different sound patterns that are created in other-language versions). The translator must project himself into the mind of the original author; he must transport himself into an entirely different world of relationships between sounds and meanings, and at the same time he must establish an equivalence between one infinitely complex system and another. Since no two languages are truly equivalent in anything except the simplest terms, this is a most difficult accomplishment. Certain writers are exceptionally difficult to translate. There are no satisfactory English versions, for example, of the Latin of Catullus, the French of Baudelaire, the Russian of Pushkin, or of the majority of Persian and Arabic poetry. The splendour of Sophocles' Greek, of Plato at his best, is barely suggested even in the finest English versions. On the other hand, the Germans insist that Shakespeare is better in German than he is in English, a humorous exaggeration perhaps. But again, Shakespeare is resistant to translation into French. His English seems to lack equivalents in that language.

The very greatest translations may become classics in their own right, of enduring literary excellence (the King James Version of the Bible, appearing in 1611, is an outstanding example), but on the whole the approximate equivalence of most translations to their originals seems to have a very short life. The original work remains the same, of lasting value to its own people, but the translation becomes out of date with each succeeding generation as the language and criteria of literary taste change. Nothing demonstrates the complexity of literary language more vividly. An analogous process takes place when a reader experiences a literary work in his own language; each generation gets a "new version" from its own classics.

Yet the values of great literature are more fundamental than complexity and subtleties of meaning arising from language alone. Works far removed from contemporary man in time and in cultural background, composed in a variety of languages utterly different from one another in structure, have nevertheless been translated successfully enough to be deeply moving. The 20th century has seen an immense mass of the oral literature of preliterate peoples and of the writings of all the great civilizations translated into modern languages. Understanding the growth of literature and its forms in other civilizations has greatly enriched the understanding of our own.

Craftsmanship. *Prosody.* Literature, like music, is an art of time, or "tempo": it takes time to read or listen to, and it usually presents events or the development of ideas or the succession of images or all these together in time. The craft of literature, indeed, can be said to be in part the manipulation of a structure in time, and so the simplest element of marking time, rhythm, is therefore of basic importance in both poetry and prose. Prosody, which is the science of versification, has for its subject the materials of poetry and is concerned almost entirely with the laws of metre, or rhythm in the narrowest sense. It deals with the patterning of sound in time; the number, length, accent, and pitch of syllables; and the modifications of rhythm by vowels and consonants. In most poetry, certain basic rhythms are repeated with modifications (that is to say, the poem rhymes or scans or both) but not in all. It most obviously does neither in the case of the "free forms" of modern poetry; but neither does it in the entire poetry of whole cultures. Since lyric poetry is either the actual text of song or else is immediately derived from song, it is regular in structure nearly everywhere in the world, although the elements of patterning that go into producing its rhythm may vary. The most important of these elements in English poetry, for example, have been accent, grouping of syllables (called feet), number of syllables in the line,

and rhyme at the end of a line (and sometimes within it). Other elements such as pitch, resonance, repetition of vowels (assonance), repetition of consonants (alliteration), and breath pauses (cadence) have also been of great importance in distinguishing successful poetry from doggerel verse, but on the whole they are not as important as the former, and poets have not always been fully conscious of their use of them. Greek and Latin poetry was consciously patterned on the length of syllables (long or short) rather than on their accent; but all the considerations of "sound" (such as assonance and alliteration) entered into the aesthetically satisfactory structure of a poem. Similarly, both the French and Japanese were content simply to count the syllables in a line—but again, they also looked to all the "sound" elements.

The rhythms of prose are more complicated, though not necessarily more complex, than those of poetry. The rules of prose patterning are less fixed; patterns evolve and shift indefinitely and are seldom repeated except for special emphasis. So the analysis of prose rhythm is more difficult to make than, at least, the superficial analysis of poetry.

Structure. The craft of writing involves more than mere rules of prosody. The work's structure must be manipulated to attract the reader. First, the literary situation has to be established. The reader must be directly related to the work, placed in it—given enough information on who, what, when, or why—so that his attention is caught and held (or, on the other hand, he must be deliberately mystified, to the same end).

Aristotle gave a formula for dramatic structure that can be generalized to apply to most literature: presentation, development, complication, crisis, and resolution. Even lyric poems can possess plot in this sense, but by no means are all literary works so structured, nor does such structure ensure their merit—it can be safely said that westerns, detective stories, and cheap melodramas are more likely to follow strictly the rules of Aristotle's *Poetics* than are great novels. Nevertheless, the scheme does provide a norm from which there is infinite variation. Neoclassical dramatists and critics, especially in 17th-century France, derived from Aristotle what they called the unities of time, action, and place. This meant that the action of a play should not spread beyond the events of one day and, best of all, should be confined within the actual time of performance. Nor should the action move about too much from place to place—best only to go from indoors to outdoors and back. There should be only one plot line, which might be relieved by a subplot, usually comic. These three unities—of time, place, and action—do not occur in Aristotle and are certainly not observed in Classical Greek tragedy. They are an invention of Renaissance critics, some of whom went even further, insisting also on what might be called a unity of mood. To this day there are those who, working on this principle, object to Shakespeare's use of comic relief within the tragic action of his plays—to the porter in *Macbeth*, for instance, or the gravediggers in *Hamlet*.

Assiduous critics have found elaborate architectural structures in quite diffuse works—including Miguel de Cervantes' *Don Quixote* (1605–15), Sterne's *Tristram Shandy* (1759–67), Casanova's *Icosameron* (1788; 1928). But their "discoveries" are too often put there after the event. Great early novels such as the Chinese *Dream of the Red Chamber* (1754; first published in English 1929) and the Japanese *Tale of Genji* (early 11th century) usually develop organically rather than according to geometrical formulas, one incident or image spinning off another. Probably the most tightly structured work, in the Neoclassicists' sense, is the Icelandic *Njál's saga*.

The 19th century was the golden age of the novel, and most of the more famous examples of the form were systematically plotted, even where the plot structure simply traced the growth in personality of an individual hero or heroine. This kind of novel, of which in their very diverse ways Stendhal's *The Red and the Black* (1830) and Dickens' *David Copperfield* (1850) are great examples, is known as *Bildungsroman*. Gustave Flaubert's *Madame Bovary* (1857) is as rigorously classicist in form as the 17th-century plays of Racine and Corneille, which were the high point of the French classical theatre, although

Flaubert obeys laws more complex than those of the Aristotelians. Novels such as Tolstoy's *War and Peace* (1865–69), Dostoyevsky's *Brothers Karamazov* (1880), and the works of Balzac owe much of their power to their ability to overwhelm the reader with a massive sense of reality. The latter 19th and early 20th centuries witnessed an attack on old forms, but what the new writers evolved was simply a new architecture. A novel like James Joyce's *Ulysses* (1922), which takes place in a day and an evening, is one of the most highly structured ever written. Novelists such as Joseph Conrad, Ford Madox Ford, Virginia Woolf, and, to a lesser extent, Henry James developed a multiple-aspect narrative, sometimes by using time shifts and flashbacks and by writing from different points of view, sometimes by using the device (dating back to Classical Greek romances) of having one or more narrators as characters within the story. (This technique, which was first perfected in the verse novels of Robert Browning, in fact reached its most extreme development in the English language in poetry: in Ezra Pound's *Cantos*, T.S. Eliot's *The Waste Land*, William Carlos Williams' *Paterson*, and the many long poems influenced by them.)

CONTENT OF LITERATURE

The word as symbol. The content of literature is as limitless as the desire of human beings to communicate with one another. The thousands of years, perhaps hundreds of thousands, since the human species first developed speech have seen built up the almost infinite systems of relationships called languages. A language is not just a collection of words in an unabridged dictionary but the individual and social possession of living human beings, an inexhaustible system of equivalents, of sounds to objects and to one another. Its most primitive elements are those words that express direct experiences of objective reality, and its most sophisticated are concepts on a high level of abstraction. Words are not only equivalent to things, they have varying degrees of equivalence to one another. A symbol, says the dictionary, is something that stands for something else or a sign used to represent something. "as the lion is the symbol of courage, the cross the symbol of Christianity." In this sense all words can be called symbols, but the examples given—the lion and the cross—are really metaphors: that is, symbols that represent a complex of other symbols, and which are generally negotiable in a given society (just as money is a symbol for goods or labour). Eventually a language comes to be, among other things, a huge sea of implicit metaphors, an endless web of interrelated symbols. As literature, especially poetry, grows more and more sophisticated, it begins to manipulate this field of suspended metaphors as a material in itself, often as an end in itself. Thus, there emerge forms of poetry (and prose, too) with endless ramifications of reference, as in Japanese waka and haiku, some ancient Irish and Norse verse, and much of the poetry written in western Europe since the time of Baudelaire that is called modernist. It might be supposed that, at its most extreme, this development would be objective, constructive—aligning it with the critical theories stemming from Aristotle's *Poetics*. On the contrary, it is romantic, subjective art, primarily because the writer handles such material instinctively and subjectively, approaches it as the "collective unconscious," to use the term of the psychologist Carl Jung, rather than with deliberate rationality.

Themes and their sources. By the time literature appears in the development of a culture, the society has already come to share a whole system of stereotypes and archetypes: major symbols standing for the fundamental realities of the human condition, including the kind of symbolic realities that are enshrined in religion and myth. Literature may use such symbols directly, but all great works of literary art are, as it were, original and unique myths. The world's great classics evoke and organize the archetypes of universal human experience. This does not mean, however, that all literature is an endless repetition of a few myths and motives, endlessly retelling the first stories of civilized man, repeating the Sumerian *Epic of Gilgamesh* or Sophocles' *Oedipus the King*. The subject matter of literature is as wide as human experience it-

Narrative devices in the novel

Unities of time, place, and action

Use of symbols and myths in literature

self. Myths, legends, and folktales lie at the beginning of literature, and their plots, situations, and allegorical (metaphorical narrative) judgments of life represent a constant source of literary inspiration that never fails. This is so because mankind is constant—people share a common physiology. Even social structures, after the development of cities, remain much alike. Whole civilizations have a life pattern that repeats itself through history. Jung's term "collective unconscious" really means that mankind is one species, with a common fund of general experience. Egyptian scribes, Japanese bureaucrats, and junior executives in New York City live and respond to life in the same ways; the lives of farmers or miners or hunters vary only within narrow limits. Love is love and death is death, for a South African Bushman and a French Surrealist alike. So the themes of literature have at once an infinite variety and an abiding constancy. They can be taken from myth, from history, or from contemporary occurrence, or they can be pure invention (but even if they are invented, they are nonetheless constructed from the constant materials of real experience, no matter how fantastic the invention).

The writer's personal involvement. As time goes on, literature tends to concern itself more and more with the interior meanings of its narrative, with problems of human personality and human relationships. Many novels are fictional, psychological biographies which tell of the slowly achieved integration of the hero's personality or of his disintegration, of the conflict between self-realization and the flow of events and the demands of other people. This can be presented explicitly, where the characters talk about what is going on in their heads, either ambiguously and with reserve, as in the novels of Henry James, or overtly, as in those of Dostoyevsky. Alternatively, it can be presented by a careful arrangement of objective facts, where psychological development is described purely in terms of behaviour and where the reader's subjective response is elicited by the minute descriptions of physical reality, as in the novels of Stendhal and the greatest Chinese novels like the *Dream of the Red Chamber*, which convince the reader that through the novel he is seeing reality itself, rather than an artfully contrived semblance of reality.

Literature, however, is not solely concerned with the concrete, with objective reality, with individual psychology, or with subjective emotion. Some deal with abstract ideas or philosophical conceptions. Much purely abstract writing is considered literature only in the widest sense of the term, and the philosophical works that are ranked as great literature are usually presented with more or less of a sensuous garment. Thus, Plato's *Dialogues* rank as great literature because the philosophical material is presented in dramatic form, as the dialectical outcome of the interchange of ideas between clearly drawn, vital personalities, and because the descriptive passages are of great lyric beauty. Karl Marx's *Das Kapital* (1867–95) approaches great literature in certain passages in which he expresses the social passion he shares with the Hebrew prophets of the Old Testament. Euclid's *Elements* and St. Thomas Aquinas' *Summa theologiae* give literary, aesthetic satisfaction to some people because of their purity of style and beauty of architectonic construction. In short, most philosophical works that rank as great literature do so because they are intensely human. The reader responds to Pascal's *Pensées*, to Montaigne's *Essays*, and to Marcus Aurelius' *Meditations* as he would to living men. Sometimes the pretense of purely abstract intellectual rigour is in fact a literary device. The writings of the 20th-century philosopher Ludwig Wittgenstein, for example, owed much of their impact to this approach, while the poetry of Paul Valéry borrows the language of philosophy and science for its rhetorical and evocative power.

Relation of form to content. Throughout literary history, many great critics have pointed out that it is artificial to make a distinction between form and content, except for purposes of analytical discussion. Form determines content. Content determines form. The issue is, indeed, usually only raised at all by those critics who are more interested in politics, religion, or ideology than in literature; thus, they object to writers who they feel sacrifice ideological orthodoxy for formal perfection, message for style.

Style. But style cannot really be said to exist on paper at all; it is the way the mind of the author expresses itself in words. Since words represent ideas, there cannot be abstract literature unless a collection of nonsense syllables can be admitted as literature. Even the most avant-garde writers associated with the Cubist or nonobjective painters used language, and language is meaning, though the meaning may be incomprehensible. Oscar Wilde and Walter Pater, the great 19th-century exponents of "art for art's sake," were in fact tireless propagandists for their views, which dominate their most flowery prose. It is true that great style depends on the perfect matching of content and form, so that the literary expression perfectly reflects the writer's intention; "poor style" reveals the inability of a writer to match the two—in other words, reveals his inability to express himself. This is why we say that "style expresses the man." The veiled style of Henry James, with its subtleties, equivocations, and qualifications, perfectly reflects his complicated and subtle mind and his abiding awareness of ambiguity in human motives. At the other extreme, the style of the early 20th-century American novelist Theodore Dreiser—bumbling, clumsy, dogged, troubled—perfectly embodies his own attitudes toward life and is, in fact, his constant judgment of his subject matter. Sometimes an author, under the impression that he is simply polishing his style, may completely alter his content. As Flaubert worked over the drafts of *Madame Bovary*, seeking always the apposite word that would precisely convey his meaning, he lifted his novel from a level of sentimental romance to make it one of the great ironic tragedies of literature. Yet, to judge from his correspondence, he seems never to have been completely aware of what he had done, of the severity of his own irony.

Literature may be an art, but writing is a craft, and a craft must be learned. Talent, special ability in the arts, may appear at an early age; the special personality called genius may indeed be born, not made. But skill in matching intention and expression comes with practice. Naïve writers, "naturals" like the 17th-century English diarist Samuel Pepys, the late 18th-century French naïf Restif de la Bretonne, the 20th-century American novelist Henry Miller, are all deservedly called stylists, although their styles are far removed from the deliberate, painstaking practice of a Flaubert or a Turgenev. They wrote spontaneously whatever came into their heads; but they wrote constantly, voluminously, and were, by their own standards, skilled practitioners.

Objective-subjective expression. There are certain forms of literature that do not permit such highly personal behaviour—for instance, formal lyric poetry and classic drama. In these cases the word "form" is used to mean a predetermined structure within whose mold the content must be fitted. These structures are, however, quite simple and so cannot be said to determine the content. Racine and Corneille were contemporaries; both were Neoclassic French dramatists; both abided by all the artificial rules—usually observing the "unities" and following the same strict rules of prosody. Yet their plays, and the poetry in which they are written, differ completely. Corneille is intellectually and emotionally a Neoclassicist—clear and hard, a true objectivist, sure of both his verse and the motivations of his characters. Racine was a great romantic long before the age of Romanticism. His characters are confused and tortured; his verse throbs like the heartbeats of his desperate heroines. He is a great sentimentalist in the best and deepest meaning of that word. His later influence on poets like Baudelaire and Paul Valéry is due to his mastery of sentimental expression, not, as they supposed, to his mastery of Neoclassic form.

Verse on any subject matter can of course be written purely according to formula. The 18th century in England saw all sorts of prose treatises cast in rhyme and metre, but this was simply applied patterning. (Works such as *The Botanic Garden* [2 vol., 1794–95] by Erasmus Darwin should be sharply distinguished from James Thomson's *The Seasons* [1726–30], which is true poetry, not versified natural history—just as Virgil's *Georgics* is not an agricultural handbook.) Neoclassicism, especially in its 18th-century developments, confused—for ordinary minds, at

Coordination of content and form in good style

Presentation of characters

Predetermined structure of classic drama

any rate—formula with form and so led to the revolt called Romanticism. The leading theorists of that revolt, the poets William Wordsworth and Samuel Taylor Coleridge, in the "Preface" (1800) to *Lyrical Ballads* urged the observance of a few simple rules basic to all great poetry and demanded a return to the integrity of expressive form. A similar revolution in taste was taking place all over Europe and also in China (where the narrow pursuit of formula had almost destroyed poetry). The Romantic taste could enjoy the "formlessness" of William Blake's prophetic books, or Walt Whitman's *Leaves of Grass*, or the loose imagination of Shelley—but careful study reveals that these writers were not formless at all. Each had his own personal form.

Form and formlessness in modern literature

Time passes and the pendulum of taste swings. In the mid-20th century, Paul Valéry, T.S. Eliot, and Yvor Winters would attack what the latter called "the fallacy of expressive form," but this is itself a fallacy. All form in literature is expressive. All expression has its own form, even when the form is a deliberate quest of formlessness. (The automatic writing cultivated by the surrealists, for instance, suffers from the excessive formalism of the unconscious mind and is far more stereotyped than the poetry of the Neoclassicist Alexander Pope.) Form simply refers to organization, and critics who attack form do not seem always to remember that a writer organizes more than words. He organizes experience. Thus, his organization stretches far back in his mental process. Form is the other face of content, the outward, visible sign of inner spiritual reality.

LITERATURE AND ITS AUDIENCE

Folk and elite literatures. In preliterate societies oral literature was widely shared; it saturated the society and was as much a part of living as food, clothing, shelter, or religion. In barbaric societies, the minstrel might be a courtier of the king or chieftain, and the poet who composed liturgies might be a priest. But the oral performance itself was accessible to the whole community. As society evolved its various social layers, or classes, an "elite" literature began to be distinguishable from the "folk" literature of the people. With the invention of writing this separation was accelerated until finally literature was being experienced individually by the elite (reading a book), while folklore and folk song were experienced orally and more or less collectively by the illiterate common people.

Elite literature continuously refreshes itself with materials drawn from the popular. Almost all poetic revivals, for instance, include in their programs a new appreciation of folk song, together with a demand for greater objectivity. On the other hand folk literature borrows themes and, very rarely, patterns from elite literature. Many of the English and Scottish ballads that date from the end of the Middle Ages and have been preserved by oral tradition share plots and even turns of phrase with written literature. A very large percentage of these ballads contain elements that are common to folk ballads from all over western Europe; central themes of folklore, indeed, are found all over the world. Whether these common elements are the result of diffusion is a matter for dispute. They do, however, represent great psychological constants, archetypes of experience common to the human species, and so these constants are used again and again by elite literature as it discovers them in folklore.

Modern popular literature. There is a marked difference between true popular literature, that of folklore and folk song, and the popular literature of modern times. Popular literature today is produced either to be read by a literate audience or to be enacted on television or in the cinema; it is produced by writers who are members, however lowly, of an elite corps of professional literates. Thus, popular literature no longer springs from the people; it is handed to them. Their role is passive. At the best they are permitted a limited selectivity as consumers.

Certain theorists once believed that folk songs and even long, narrative ballads were produced collectively, as has been said in mockery "by the tribe sitting around the fire and grunting in unison." This idea is very much out of date. Folk songs and folk tales began somewhere in one

human mind. They were developed and shaped into the forms in which they are now found by hundreds of other minds as they were passed down through the centuries. Only in this sense were they "collectively" produced. During the 20th century, folklore and folk speech have had a great influence on elite literature—on writers as different as Franz Kafka and Carl Sandburg, Selma Lagerlöf and Kawabata Yasunari, Martin Buber and Isaac Bashevis Singer. Folk song has always been popular with bohemian intellectuals, especially political radicals (who certainly are an elite). Since World War II the influence of folk song upon popular song has not just been great; it has been determinative. Almost all "hit" songs since the mid-century have been imitation folk songs; and some authentic folk singers attract immense audiences.

Influence of folklore on modern literature

Popular fiction and drama, westerns and detective stories, films and television serials, all deal with the same great archetypal themes as folktales and ballads, though this is seldom due to direct influence; these are simply the limits within which the human mind works. The number of people who have elevated the formulas of popular fiction to a higher literary level is surprisingly small. Examples are H.G. Wells's early science fiction, the western stories of Gordon Young and Ernest Haycox, the detective stories of Sir Arthur Conan Doyle, Georges Simenon, and Raymond Chandler.

The latter half of the 20th century has seen an even greater change in popular literature. Writing is a static medium: that is to say, a book is read by one person at a time; it permits recollection and anticipation; the reader can go back to check a point or move ahead to find out how the story ends. In radio, television, and the cinema the medium is fluent; the audience is a collectivity and is at the mercy of time. It cannot pause to reflect or to understand more fully without missing another part of the action, nor can it go back or forward. Marshall McLuhan in his book *Understanding Media* (1964) became famous for erecting a whole structure of aesthetic, sociological, and philosophical theory upon this fact. But it remains to be seen whether the new, fluent materials of communication are going to make so very many changes in civilization, let alone in the human mind—mankind has, after all, been influenced for thousands of years by the popular, fluent arts of music and drama. Even the most transitory television serial was written down before it was performed, and the script can be consulted in the files. Before the invention of writing, all literature was fluent because it was contained in people's memory. In a sense it was more fluent than music, because it was harder to remember. Man in mass society becomes increasingly a creature of the moment, but the reasons for this are undoubtedly more fundamental than his forms of entertainment.

LITERATURE AND ITS ENVIRONMENT

Social and economic conditions. Literature, like all other human activities, necessarily reflects current social and economic conditions. Class stratification was reflected in literature as soon as it had appeared in life. Among the American Indians, for instance, the chants of the shaman, or medicine man, differ from the secret, personal songs of the individual, and these likewise differ from the group songs of ritual or entertainment sung in community. In the Heroic Age, the epic tales of kings and chiefs that were sung or told in their barbaric courts differed from the folktales that were told in peasant cottages.

Reflection of class distinction in literature

The more cohesive a society, the more the elements—and even attitudes—evolved in the different class strata are interchangeable at all levels. In the tight clan organization that existed in late medieval times at the Scottish border, for example, heroic ballads telling of the deeds of lords and ladies were preserved in the songs of the common people. But where class divisions are unbridgeable, elite literature is liable to be totally separated from popular culture. An extreme example is the classic literature of the Roman Empire. Its forms and its sources were largely Greek—it even adopted its laws of verse patterning from Greek models, even though these were antagonistic to the natural patterns of the Latin language—and most of the sophisticated works of the major Latin authors were con-

pletely closed to the overwhelming majority of people of the Roman Empire.

Printing has made all the difference in the negotiability of ideas. The writings of the 18th-century French writers Voltaire, Rousseau, and Diderot were produced from and for almost as narrow a caste as the Roman elite, but they were printed. Within a generation they had penetrated the entire society and were of vital importance in revolutionizing it.

Class distinctions in the literature of modern times exist more in the works themselves than in their audience. Although Henry James wrote about the upper classes and Émile Zola about workingmen, both were, in fact, members of an elite and were read by members of an elite—moreover, in their day, those who read Zola certainly considered themselves more of an elite than did the readers of Henry James. The ordinary people, if they read at all, preferred sentimental romances and “penny dreadfuls.” Popular literature had already become commercially produced entertainment literature, a type which today is also provided by television scripts.

The elite who read serious literature are not necessarily members of a social or economic upper class. It has been said of the most ethereal French poet, Stéphane Mallarmé, that in every French small town there was a youth who carried his poems in his heart. These poems are perhaps the most “elite” product of western European civilization, but the “youths” referred to were hardly the sons of dukes or millionaires. (It is a curious phenomenon that, since the middle of the 18th century in Europe and in the United States, the majority of readers of serious literature—as well as of entertainment literature—have been women. The extent of the influence that this audience has exerted on literature itself must be immense.)

National and group literature. Hippolyte Taine, the 19th-century French critic, evolved an ecological theory of literature. He looked first and foremost to the national characteristics of western European literatures, and he found the source of these characteristics in the climate and soil of each respective nation. His *History of English Literature* (5 vol., 1863–69) is an extensive elaboration of these ideas. It is doubtful that anyone today would agree with the simplistic terms in which Taine states his thesis. It is obvious that Russian literature differs from English or French from German. English books are written by Englishmen, their scenes are commonly laid in England, they are usually about Englishmen and they are designed to be read by Englishmen—at least in the first instance. But modern civilization becomes more and more a world civilization, wherein works of all peoples flow into a general fund of literature. It is not unusual to read a novel by a Japanese author one week and one by a black writer from West Africa the next. Writers are themselves affected by this cross-fertilization. Certainly, the work of the great 19th-century Russian novelists has had more influence on 20th-century American writers than has the work of their own literary ancestors. Poetry does not circulate so readily, because catching its true significance in translation is so very difficult to accomplish. Nevertheless, for the past 100 years or so, the influence of French poetry upon all the literatures of the civilized world has not just been important, it has been preeminent. The tendentious elements of literature—propaganda for race, nation, or religion—have been more and more eroded in this process of wholesale cultural exchange.

Popular literature on the other hand is habitually tendentious both deliberately and unconsciously. It reflects and stimulates the prejudices and parochialism of its audience. Most of the literary conflicts that have seized the totalitarian countries during the 20th century stem directly from relentless efforts by the state to reduce elite literature to the level of the popular. The great proletarian novels of our time have been produced, not by Russians, but by American blacks, Japanese, Germans, and—most proletarian of all—a German-American living in Mexico, B. Traven. Government control and censorship can inhibit literary development, perhaps deform it a little, and can destroy authors outright; but, whether in the France of Louis XIV or in the Soviet Union of the 20th century,

it cannot be said to have a fundamental effect upon the course of literature.

The writer's position in society. A distinguishing characteristic of modern literature is the peculiar elite which it has itself evolved. In earlier cultures the artist, though he may have been neurotic at times, thought of himself as part of his society and shared its values and attitudes. Usually the clerkly caste played a personal, important role in society. In the modern industrial civilization, however, “scribes” became simply a category of skilled hired hands. The writer shared few of the values of the merchant or the entrepreneur or manager. And so the literary and artistic world came to have a subculture of its own. The antagonism between the two resultant sets of values is the source of what we call alienation—among the intellectuals at least (the alienation of the common man in urban, industrial civilization from his work, from himself, and from his fellows is another matter, although its results are reflected and intensified in the alienation of the elite). For about 200 years now, the artistic environment of the writer has not usually been shared with the general populace. The subculture known as bohemia and the literary and artistic movements generated in its little special society have often been more important—at least in the minds of many writers—than the historical, social, and economic movements of the culture as a whole. Even massive historical change is translated into these terms—the Russian Revolution, for instance, into Communist-Futurism, Constructivism, Socialist Realism. Western European literature could be viewed as a parade of movements—Romanticism, Realism, Naturalism, Futurism, Structuralism, and so on indefinitely. Some of the more journalistic critics, indeed, have delighted to regard it in such a way. But after the manifestos have been swept away, the meetings adjourned, the literary cafés of the moment lost their popularity, the turmoil is seen not to have made so very much difference. The Romantic Théophile Gautier (1811–72) and the Naturalist Émile Zola (1840–1902) have more in common than they have differences, and their differences are rather because of changes in society as a whole than because of conflicting literary principles.

At first, changes in literary values are appreciated only at the upper levels of the literary elite itself, but often, within a generation, works once thought esoteric are being taught as part of a school syllabus. Most cultivated people once thought James Joyce's *Ulysses* incomprehensible or, where it was not, obscene. Today his methods and subject matter are commonplace in the commercial fiction of the mass culture. A few writers remain confined to the elite. Mallarmé is a good example—but he would have been just as ethereal had he written in the simplest French of direct communication. His subtleties are ultimately grounded in his personality.

Literature and the other arts. Literature has an obvious kinship with the other arts. Presented, a play is drama; read, a play is literature. Most important films have been based upon written literature, usually novels, although all the great epics and most of the great plays have been filmed at some time and thus have stimulated the younger medium's growth. Conversely, the techniques required in writing for film have influenced many writers in structuring their novels and have affected their style. Most popular fiction is written with “movie rights” in mind, and these are certainly a consideration with most modern publishers. Literature provides the libretto for operas, the theme for tone poems—even so anomalous a form as Nietzsche's *Thus Spake Zarathustra* was interpreted in music by Richard Strauss—and of course it provides the lyrics of songs. Many ballets and modern dances are based on stories or poems. Sometimes, music and dance are accompanied by a text read by a speaker or chanted by a chorus. The mid-19th century was the heyday of literary, historical, and anecdotal painting, though, aside from the Surrealists, this sort of thing died out in the 20th century. Cross-fertilization of literature and the arts now takes place more subtly, mostly in the use of parallel techniques—the rational dissociation of the Cubists or the spontaneous action painting of the Abstract Expressionists, for example, which flourished at the same time as

Literary
subcultures
and
movements

Hippolyte
Taine's
ecological
theory of
literature

Effect of
censorship
on literary
develop-
ment

the free-flowing uncorrected narratives of some novelists in the 1950s and '60s.

LITERATURE AS A COLLECTION OF GENRES

Critics have invented a variety of systems for treating literature as a collection of genres. Often these genres are artificial, invented after the fact with the aim of making literature less sprawling, more tidy. Theories of literature must be based upon direct experience of the living texts and so be flexible enough to contain their individuality and variety. Perhaps the best approach is historical, or genetic. What actually happened, and in what way did literature evolve up to the present day?

Oral
literature

There is a surprising variety of oral literature among surviving preliterate peoples, and, as the written word emerges in history, the indications are that the important literary genres all existed at the beginning of civilized societies: heroic epic; songs in praise of priests and kings; stories of mystery and the supernatural; love lyrics; personal songs (the result of intense meditation); love stories; tales of adventure and heroism (of common peoples, as distinct from the heroic epics of the upper classes); satire (which was dreaded by barbaric chieftains); satirical combats (in which two poets or two personifications abused one another and praised themselves); ballads and folktales of tragedy and murder; folk stories, such as the tale of the clever boy who performs impossible tasks, outwits all his adversaries, and usually wins the hand of the king's daughter; animal fables like those attributed to Aesop (the special delight of Black Africa and Indian America); riddles, proverbs, and philosophical observations; hymns, incantations, and mysterious songs of priests; and finally actual mythology—stories of the origin of the world and the human race, of the great dead, and of the gods and demigods.

Epic. The true heroic epic never evolved far from its preliterate origins, and it arose only in the Heroic Age which preceded a settled civilization. The conditions reflected in, say, the *Iliad* and *Odyssey* are much the same as those of the Anglo-Saxon *Beowulf*, the German *Nibelungenlied*, or the Irish stories of Cú Chulainn. The literary epic is another matter altogether. Virgil's *Aeneid*, for instance, or Milton's *Paradise Lost* are products of highly sophisticated literary cultures. Many long poems sometimes classified as epic literature are no such thing—Dante's *La divina commedia* (*The Divine Comedy*), for example, is a long theological, philosophical, political, moral, and mystical poem. Dante considered it to be a kind of drama which obeyed the rules of Aristotle's *Poetics*. Goethe's *Faust* is in dramatic form and is sometimes even staged—but it is really a philosophical poetic novel. Modern critics have described long poems such as T.S. Eliot's *Waste Land* and Ezra Pound's *Cantos* as "philosophical epics." There is nothing epic about them; they are reveries, more or less philosophical.

Lyric poetry. Lyric poetry never gets far from its origins, except that some of its finest examples—Medieval Latin, Provençal, Middle High German, Middle French, Renaissance—which today are only read, were actually written to be sung. In the 20th century, however, popular songs of great literary merit have become increasingly common—for example, the songs of Bertolt Brecht and Kurt Weill in German, of Georges Brassens and Anne Sylvestre in French, and of Leonard Cohen, Bob Dylan, and Joni Mitchell. It is interesting to note that, in periods when the culture values artificiality, the lyric becomes stereotyped. Then, after a while, the poets revolt and, usually turning to folk origins, restore to lyric poetry at least the appearance of naturalness and spontaneity.

Diverse
forms of
satire

Satire. The forms of satire are as manifold as those of literature itself—from those of the mock epic to the biting epigram. A great many social and political novels of today would have been regarded as satire by the ancients. Many of the great works of all time are satires, but in each case they have risen far above their immediate satirical objectives. The 16th-century medieval satire on civilization, the *Gargantua and Pantagruel* of Rabelais, grew under the hand of its author into a great archetypal myth of the lust for life. Cervantes' *Don Quixote*, often called the greatest work of prose fiction in the West, is superficially a satire

of the sentimental romance of knightly adventure. But, again, it is an archetypal myth, telling the adventures of the soul of man—of the individual—in the long struggle with what is called the human condition. *The Tale of Genji* by Murasaki Shikibu has sometimes been considered by obtuse critics as no more than a satire on the sexual promiscuity of the Heian court. In fact, it is a profoundly philosophical, religious, and mystical novel.

Prose fiction. Extended prose fiction is the latest of the literary forms to develop. We have romances from classical Greek times that are as long as short novels; but they are really tales of adventure—vastly extended anecdotes. The first prose fiction of any psychological depth is the *Satyricon*, almost certainly attributed to Petronius Arbiter (died AD 65/66). Though it survives only in fragments, supposedly one-eleventh of the whole, even these would indicate that it is one of the greatest picaresque novels, composed of loosely connected episodes of robust and often erotic adventure. The other great surviving fiction of classical times is the *Metamorphoses* (known as *The Golden Ass*) by Apuleius (2nd century AD). In addition to being a picaresque adventure story, it is a criticism of Roman society, a celebration of the religion of Isis, and an allegory of the progress of the soul. It contains the justly celebrated story of Cupid and Psyche, a myth retold with psychological subtlety. Style has much to do with the value and hence the survival of these two works. They are written in prose of extraordinary beauty, although it is by no means of "classical" purity. The prose romances of the Middle Ages are closely related to earlier heroic literature. Some, like Sir Thomas Malory's 15th-century *Le Morte Darthur*, are retellings of heroic legend in terms of the romantic chivalry of the early Renaissance, a combination of barbaric, medieval, and Renaissance sensibility which, in the tales of Tristram and Iseult and Lancelot and Guinevere, produced something not unlike modern novels of tragic love.

The earliest
prose
fiction

The Western novel is a product of modern civilization, although in the Far East novels began a separate development as early as the 10th century. Extended prose works of complex interpersonal relations and motivations begin in 17th-century France with *The Princess of Cleves* (1678) by Madame de La Fayette. Eighteenth-century France produced an immense number of novels dealing with love analysis but none to compare with Madame de La Fayette's until Pierre Choderlos de Laclos wrote *Les Liaisons dangereuses* (1782). This was, in form, an exchange of letters between two corrupters of youth; but, in intent, it was a savage satire of the *ancien régime* and a heart-rending psychological study. The English novel of the 18th century was less subtle, more robust—vulgar in the best sense—and is exemplified by Henry Fielding's *Tom Jones* (1749) and Laurence Sterne's *Tristram Shandy*. The 19th century was the golden age of the novel. It became ever more profound, complex, and subtle (or, on the other hand, more popular, eventful, and sentimental). By the beginning of the 20th century it had become the most common form of thoughtful reading matter and had replaced, for most educated people, religious, philosophical, and scientific works as a medium for the interpretation of life. By the late 1920s the novel had begun to show signs of decay as a form, and no works have since been produced to compare with the recent past. This may prove to be a temporarily barren period, or else the novel may be losing its energy as a narrative art form and in this sense giving way to the medium of film.

Current
state of the
novel

Drama. Like lyric poetry, drama has been an exceptionally stable literary form. Given a little leeway, most plays written by the beginning of the 20th century could be adjusted to the rules of Aristotle's *Poetics*. Before World War I, however, all traditional art forms, led by painting, began to disintegrate, and new forms evolved to take their place. In drama the most radical innovator was August Strindberg (1849–1912), and from that day to this, drama (forced to compete with the cinema) has become ever more experimental, constantly striving for new methods, materials, and, especially, ways to establish a close relationship with the audience. All this activity has profoundly modified drama as literature.

Future developments. In the 20th century the methods of poetry have also changed drastically, although the "innovator" here might be said to have been Baudelaire (1821–67). The disassociation and recombination of ideas of the Cubists, the free association of ideas of the Surrealists, dreams, trance states, the poetry of preliterate people—all have been absorbed into the practice of modern poetry. This proliferation of form is not likely to end. Effort that once was applied to perfecting a single pattern in a single form may in the future be more and more directed toward the elaboration of entirely new "multimedia" forms, employing the resources of all the established arts. At the same time, writers may prefer to simplify and polish the forms of the past with a rigorous, Neoclassicist discipline. In a worldwide urban civilization, which has taken to itself the styles and discoveries of all cultures past and present, the future of literature is quite impossible to determine.

WRITINGS ON LITERATURE

Scholarly research. Research by scholars into the literary past began almost as soon as literature itself—as soon as the documents accumulated—and for many centuries it represents almost all the scholarship that has survived. The most extensive text of the Sumerian *Epic of Gilgamesh*, the first of the world's great classics, is a late Assyrian synthesis that must have required an immense amount of research into clay tablets, written in several languages going back to the beginning of Mesopotamian civilization. Many Egyptian poems and the philosophic creation myth known as the "Memphite Theology" survive in very late texts that carefully reproduce the original language of the first dynasties. Once the function of the scribe was established as essential, he invented literary scholarship, both to secure his position and to occupy his leisure. The great epoch of literary scholarship in ancient times centred on the library (and university) of Alexandria from its foundation in 324 BC to its destruction by the Arabs in AD 640. Hellenistic Greek scholars there developed such an academic and pedantic approach to literary scholarship and scholarly literature that the term Alexandrine remains pejorative to this day. To them, however, is owed the survival of the texts of most of the Greek classics. Roman literary scholarship was rhetorical rather than analytic. With the coming of Islām, there was established across the whole warm temperate zone of the Old World a far-flung community of scholars who were at home in learned circles from India to Spain. Judaism, like Islām, was a religion of the book and of written tradition, so literary scholarship played a central role in each. The same is true of India, China, and later Japan; for sheer bulk, as well as for subtlety and insight, Oriental scholarship has never been surpassed. In a sense, the Renaissance in Europe was a cultural revolution led by literary scholars who

Ancient
literary
scholarship

discovered, revived, and made relevant again the literary heritage of Greece and Rome. In the 19th century, literary scholarship was dominated by the exhaustive, painstaking German academician, and that Germanic tradition passed to the universities of the United States. The demand that every teacher should write a master's thesis, a doctor's dissertation, and, for the rest of his career, publish with reasonable frequency learned articles and scholarly books, has led to a mass of scholarship of widely varying standards and value. Some is trivial and absurd, but the best has perfected the texts and thoroughly illuminated the significance of nearly all the world's great literature.

Literary criticism. Literary criticism, as distinguished from scholarly research, is usually itself considered a form of literature. Some people find great critics as entertaining and stimulating as great poets, and theoretical treatises of literary aesthetics can be as exciting as novels. Aristotle, Longinus, and the Roman rhetorician and critic Quintilian are still read, although Renaissance critics like the once all-powerful Josephus Scaliger are forgotten by all but specialized scholars. Later critics, such as Poe, Sainte-Beuve, Taine, Vissarion Belinsky, Matthew Arnold, Walter Bagehot, Walter Pater, and George Saintsbury, are probably read more for themselves than for their literary judgments and for their general theorizing rather than for their applications (in the case of the first three, for instance, time has confounded almost all the evaluations they made of their contemporaries). The English critics have survived because they largely confined themselves to acknowledged masterpieces and general ideas. Perhaps literary criticism can really be read as a form of autobiography. Aestheticians of literature like I.A. Richards, Sir C.M. Bowra, Paul Valéry, Suzanne Langer, and Ernst Cassirer have had an influence beyond the narrow confines of literary scholarship and have played in our time something approaching the role of general philosophers. This has been true on the popular level as well. The Dane Georg Brandes, the Americans James Gibbons Huneker, H.L. Mencken, and Edmund Wilson—these men have been social forces in their day. Literary criticism can play its role in social change. In Japan, the overthrow of the shogunate, the restoration of the emperor, and the profound change in the Japanese social sensibility begins with the literary criticism of Motoori Norinaga (1730–1801). The 19th-century revolution in theology resulted from the convergence of Darwinian theories of evolution and the technical and historical criticism of the Bible that scholars had undertaken. For many modern intellectuals, the literary quarterlies and weeklies, with their tireless discussions of the spiritual significance and formal characteristics of everything from the greatest masterpiece to the most ephemeral current production, can be said to have filled the place of religion, both as rite and dogma. (K.Re.)

POETRY

The nature of poetry

Poetry is a vast subject, as old as history and older, present wherever religion is present, possibly—under some definitions—the primal and primary form of languages themselves. The present section means only to describe in as general a way as possible certain properties of poetry and of poetic thought regarded as in some sense independent modes of the mind. Naturally, not every tradition nor every local or individual variation can be—or need be—included, but the article illustrates by examples of poetry ranging between nursery rhyme and epic. (For particular information about various types of poetry, see below *Prosody*, *Ballad*, and *Epic*.)

ATTEMPTS TO DEFINE POETRY

Poetry is the other way of using language. Perhaps in some hypothetical beginning of things it was the only way of using language or simply was language *tout court*, prose being the derivative and younger rival. Both poetry and language are fashionably thought to have belonged to rit-

ual in early agricultural societies; and poetry in particular, it has been claimed, arose at first in the form of magical spells recited to ensure a good harvest. Whatever the truth of this hypothesis, it blurs a useful distinction: by the time there begins to be a separate class of objects called poems, recognizable as such, these objects are no longer much regarded for their possible yam-growing properties, and such magic as they may be thought capable of has retired to do its business upon the human spirit and not directly upon the natural world outside.

Formally, poetry is recognizable by its greater dependence on at least one more parameter, the *line*, than appears in prose composition. This changes its appearance on the page; and it seems clear that people take their cue from this changed appearance, reading poetry aloud in a very different voice from their habitual voice, possibly because, as Ben Jonson said, poetry "speaketh somewhat above a mortal mouth." If, as a test of this description, people are shown poems printed as prose, it most often turns out that they will read the result as prose simply because it looks that way; which is to say that they are no longer

guided in their reading by the balance and shift of the line in relation to the breath as well as the syntax.

That is a minimal definition but perhaps not altogether uninformative. It may be all that ought to be attempted in the way of a definition: Poetry is the way it is because it looks that way, and it looks that way because it sounds that way and vice versa.

POETRY AND PROSE

People's reason for wanting a definition is to take care of the borderline case, and this is what a definition, as if by definition, will not do. That is, if a man asks for a definition of poetry, it will most certainly not be the case that he has never seen one of the objects called poems that are said to embody poetry; on the contrary, he is already tolerably certain what poetry in the main is, and his reason for wanting a definition is either that his certainty has been challenged by someone else or that he wants to take care of a possible or seeming exception to it: hence the perennial squabble about distinguishing poetry from prose, which is rather like distinguishing rain from snow—everyone is reasonably capable of doing so, and yet there are some weathers that are either-neither.

Similar things have been said on the question. The poet T.S. Eliot suggested that part of the difficulty lies in the fact that there is the technical term "verse" to go with the term "poetry," while there is no equivalent technical term to distinguish the mechanical part of prose and make the relation symmetrical. The French poet Paul Valéry said that prose was walking, poetry dancing. Indeed, the original two terms, *prosus* and *versus*, meant, respectively, "going straight forth" and "returning"; and that distinction does point up the tendency of poetry to incremental repetition, variation, and the treatment of many matters and different themes in a single recurrent form such as couplet or stanza.

Robert Frost said shrewdly that poetry was what got left behind in translation, which suggests a criterion of almost scientific refinement: when in doubt, translate; whatever comes through is prose, the remainder is poetry. And yet to even so acute a definition the obvious exception is a startling and a formidable one: some of the greatest poetry in the world is in the Authorized Version of the Bible, which is not only a translation but also, as to its appearance in print, identifiable neither with verse nor with prose in English but rather with a cadence owing something to both.

There may be a better way of putting the question by the simple test alluded to above. When people are presented with a series of passages drawn indifferently from poems and stories but all printed as prose, they will show a dominant inclination to identify everything they possibly can as prose. This will be true, surprisingly enough, even if the poem rhymes and will often be true even if the poem in its original typographical arrangement would have been familiar to them. The reason seems to be absurdly plain: the reader recognizes poems by their appearance on the page, and he responds to the convention whereby he recognizes them by reading them aloud in a quite different tone of voice from that which he applies to prose (which, indeed, he scarcely reads aloud at all). It should be added that he makes this distinction also without reading aloud; even in silence he confers upon a piece of poetry an attention that differs from what he gives to prose in two ways especially: in tone and in pace.

In place of further worrying over definitions, it may be both a relief and an illumination to exhibit certain plain and mighty differences between prose and poetry by a comparison. In the following passages a prose writer and a poet are talking about the same subject, growing older.

Between the ages of 30 and 90, the weight of our muscles falls by 30 percent and the power we can exert likewise The number of nerve fibres in a nerve trunk falls by a quarter. The weight of our brains falls from an average of 3.03 lb. to 2.27 lb. as cells die and are not replaced. . . . (Gordon Rattray Taylor, *The Biological Time Bomb*, 1968.)

Let me disclose the gifts reserved for age
To set a crown upon your lifetime's effort.
First, the cold friction of expiring sense

Without enchantment, offering no promise
But bitter tastelessness of shadow fruit
As body and soul begin to fall asunder.
Second, the conscious impotence of rage
At human folly, and the laceration
Of laughter at what ceases to amuse.
And last, the rending pain of re-enactment
Of all that you have done, and been

(T.S. Eliot, *Four Quartets*.)

Before objecting that a simple comparison cannot possibly cover all the possible ranges of poetry and prose compared, the reader should consider for a moment what differences are exhibited. The passages are oddly parallel, hence comparable, even in a formal sense; for both consist of the several items of a catalog under the general title of growing old. The significant differences are of tone, pace, and object of attention. If the prose passage interests itself in the neutral, material, measurable properties of the process, while the poetry interests itself in what the process will signify to someone going through it, that is not accidental but of the essence; if one reads the prose passage with an interest in being informed, noting the parallel constructions without being affected by them either in tone or in pace, while reading the poetry with a sense of considerable gravity and solemnity, that too is of the essence. One might say as tersely as possible that the difference between prose and poetry is most strikingly shown in the two uses of the verb "to fall":

The number of nerve fibres in a nerve trunk falls by a quarter
As body and soul begin to fall asunder

It should be specified here that the important differences exhibited by the comparison belong to the present age. In each period, speaking for poetry in English at any rate, the dividing line will be seen to come at a different place. In Elizabethan times the diction of prose was much closer to that of poetry than it later became, and in the 18th century authors saw nothing strange about writing in couplets about subjects that later would automatically and compulsorily belong to prose—for example, horticulture, botany, even dentistry. Here is not the place for entering into a discussion of so rich a chapter in the history of ideas; but it should be remarked that the changes involved in the relation between poetry and prose are powerfully influenced by the immense growth of science, commerce, and number in man's ways of describing, even of viewing, the world.

Returning to the comparison, it is observable that though the diction of the poem is well within what could be commanded by a moderately well-educated speaker, it is at the same time well outside the range of terms in fact employed by such a speaker in his daily occasions: it is a diction very conscious, as it were, of its power of choosing terms with an effect of peculiar precision and of combining the terms into phrases with the same effect of peculiar precision and also of combining sounds with the same effect of peculiar precision. Doubtless the precision of the prose passage is greater in the more obvious property of dealing in the measurable; but the poet attempts a precision with respect to what is not in the same sense measurable nor even in the same sense accessible to observation: the distinction is perhaps just that made by the French scientist and philosopher Blaise Pascal in discriminating the spirits of geometry and finesse; and if one speaks of "effects of precision" rather than of precision itself, that serves to distinguish one's sense that the art work is always somewhat removed from what people are pleased to call the real world, operating instead, in Immanuel Kant's shrewd formula, by exhibiting "purposefulness without purpose." To much the same point is what Samuel Taylor Coleridge remembers having learned from his schoolmaster:

I learnt from him, that Poetry, even that of the loftiest and, seemingly, that of the wildest odes, had a logic of its own, as severe as that of science; and more difficult, because more subtle, more complex, and dependent on more, and more fugitive causes. In the truly great poets, he would say, there is a reason assignable, not only for every word, but for the position of every word. (*Biographia Literaria*, ch. 1.)

Perhaps this is a somewhat exaggerated, as it is almost

Differences of tone, pace, and object of attention

Precision of poetic diction

Distin-
guishing
prose from
poetry

always an unprovable, claim, illustrating also a propensity for competing with the prestige of science on something like its own terms—but the last remark in particular illuminates the same author's terser formulation: "prose = words in the best order, poetry = the best words in the best order." This attempt at definition, impeccable because uninformative, was derived from Jonathan Swift, who had said, also impeccably and uninformatively, that style in writing was "the best words in the best order." Which may be much to the same effect as Louis Armstrong's saying, on being asked to define jazz, "Baby, if you got to ask the question, you're never going to know the answer." Or the painter Marcel Duchamp's elegant remark on what psychologists call "the problem of perception": "If no solution, then maybe no problem?" This species of gnomic, riddling remark may be determinate for the artistic attitude toward definition of every sort; and its skepticism is not confined to definitions of poetry but extends to definitions of anything whatever, directing one not to dictionaries but to experience and, above all, to use: "Anyone with a watch can tell you what time it is," said Valéry, "but who can tell you what is time?"

Happily, if poetry is almost impossible to define, it is extremely easy to recognize in experience; even untutored children are rarely in doubt about it when it appears:

Little Jack Jingle,
He used to live single,
But when he got tired of this kind of life,
He left off being single, and liv'd with his wife.

It might be objected that this little verse is not of sufficient import and weight to serve as an exemplar for poetry. It ought to be remembered, though, that it has given people pleasure so that they continued to say it until and after it was written down, nearly two centuries ago. The verse has survived, and its survival has something to do with pleasure, with delight; and while it still lives, how many more imposing works of language—epic poems, books of science, philosophy, theology—have gone down, deservedly or not, into dust and silence. It has, obviously, a form, an arrangement of sounds in relation to thoughts that somehow makes its agreeable nonsense closed, complete, and decisive. But this somewhat muddled matter of form deserves a heading and an instance all to itself.

FORM IN POETRY

People nowadays who speak of form in poetry almost always mean such externals as regular measure and rhyme, and most often they mean to get rid of these in favour of the freedom they suppose must follow upon the absence of form in this limited sense. But in fact a poem having only one form would be of doubtful interest even if it could exist. In this connection, the poet J.V. Cunningham speaks of "a convergence of forms, and forms of disparate orders," adding: "It is the coincidence of forms that locks in the poem." For a poem is composed of internal and intellectual forms as well as forms externally imposed and preexisting any particular instance, and these may be sufficient without regular measure and rhyme; if the intellectual forms are absent, as in greeting-card verse and advertising jingles, no amount of thumping and banging will supply the want.

Form, in effect, is like the doughnut that may be said to be nothing in a circle of something or something around nothing; it is either the outside of an inside, as when people speak of "good form" or "bourgeois formalism," or the inside of an outside, as in the scholastic saying that "the soul is the form of the body." Taking this principle, together with what Cunningham says of the matter, one may now look at a very short and very powerful poem with a view to distinguishing the forms, or schemes, of which it is made. It was written by Rudyard Kipling—a great poet at present somewhat sunken in reputation, probably on account of misinterpretations having to do more with his imputed politics than with his poetry—and its subject, one of a series of epitaphs for the dead of World War I, is a soldier shot by his comrades for cowardice in battle.

I could not look on Death, which being known,
Men led me to him, blindfold and alone.

The aim of the following observations and reflections is to distinguish as clearly as possible—distinguish without dividing—the feelings evoked by the subject, so grim, horrifying, tending to helpless sorrow and despair, from the feelings, which might better be thought of as meanings, evoked by careful contemplation of the poem in its manifold and somewhat subtle ways of handling the subject, leading the reader on to a view of the strange delight intrinsic to art, whose mirroring and shielding power allows him to contemplate the world's horrible realities without being turned to stone.

There is, first, the obvious external form of a rhymed, closed couplet in iambic pentameter (that is, five poetic "feet," each consisting of an unstressed followed by a stressed syllable, per line). There is, second, the obvious external form of a single sentence balanced in four grammatical units with and in counterpoint with the metrical form. There is, third, the conventional form belonging to the epitaph and reflecting back to antiquity; it is terse enough to be cut in stone and tight-lipped also, perhaps for other reasons, such as the speaker's shame. There is, fourth, the fictional form belonging to the epitaph, according to which the dead man is supposed to be saying the words himself. There is, fifth, especially poignant in this instance, the real form behind or within the fictional one, for the reader is aware that in reality it is not the dead man speaking, nor are his feelings the only ones the reader is receiving, but that the comrades who were forced to execute him may themselves have made up these two lines with their incalculably complex and exquisite balance of scorn, awe, guilt, and consideration even to tenderness for the dead soldier. There is, sixth, the metaphorical form, with its many resonances ranging from the tragic through the pathetic to irony and apology: dying in battle is spoken of in language relating it to a social occasion in drawing room or court; the coward's fear is implicitly represented as merely the timorousness and embarrassment one might feel about being introduced to a somewhat superior and majestic person, so that the soldiers responsible for killing him are seen as sympathetically helping him through a difficult moment in the realm of manners. In addition, there is, seventh, a linguistic or syntactical form, with at least a couple of tricks to it: the second clause, with its reminiscence of Latin construction, participates in the meaning by conferring a Roman stoicism and archaic gravity on the saying; remembering that the soldiers in the poem had been British schoolboys not long before, the reader might hear the remote resonance of a whole lost world built upon Greek and Roman models; and the last epithets, "blindfold and alone," while in the literal acceptance they clearly refer to the coward, show a distinct tendency to waver over and apply mysteriously to Death as well, sitting there waiting "blindfold and alone." One might add another form, the eighth, composed of the balance of sounds, from the obvious likeness in the rhyme down to subtleties and refinements beneath the ability of coarse analysis to discriminate. And even there one would not be quite at an end; an overall principle remains, the compression of what might have been epic or five-act tragedy into two lines, or the poet's precise election of a single instant to carry what the novelist, if he did his business properly, would have been hundreds of pages arriving at.

It is not at all to be inferred that the poet composed his poem in the manner of the above laborious analysis of its strands; the whole insistence, rather, is that he did not catalog eight or 10 forms and assemble them into a poem; more likely it "just came to him." But the example may serve to indicate how many modes of the mind go together in this articulation of an implied drama and the tension among many possible sentiments that might arise in response to it.

In this way, by the coincidence of forms that locks in the poem, one may see how to answer a question that often arises about poems: though their thoughts are commonplace, they themselves mysteriously are not. One may answer on the basis of the example and the inferences produced from it that a poem is not so much a thought as it is a mind: talk with it, and it will talk back, telling you

Kinds of poetic form

The coincidence of poetic forms

Experiencing poetry

many things that you might have thought for yourself but somehow didn't until it brought them together. Doubtless a poem is a much simplified model for the mind. But it might still be one of the best man has available. On this great theme, however, it will be best to proceed not by definition but by parable and interpretation.

POETRY AS A MODE OF THOUGHT:
THE PROTEAN ENCOUNTER

In the fourth book of the *Odyssey* Homer tells the following strange tale. After the war at Troy, Menelaus wanted very much to get home but was held up in Egypt for want of a wind because, as he later told Telemachus, he had not sacrificed enough to the gods. "Ever jealous the Gods are," he said, "that we men mind their dues." But because the gods work both ways, it was on the advice of a goddess, Eidothea, that Menelaus went to consult Proteus, the old one of the sea, as one might consult a travel agency.

Proteus was not easy to consult. He was herding seals, and the seals stank even through the ambrosia Eidothea had provided. And when Menelaus crept up close, disguised as a seal, and grabbed him, Proteus turned into a lion, a dragon, a leopard, a boar, a film of water, and a high-branched tree. But Menelaus managed to hang on until Proteus gave up and was himself again: whereupon Menelaus asked him the one great question: How do I get home? And Proteus told him: You had better go back to Egypt and sacrifice to the gods some more.

This story may be taken as a parable about poetry. A man has an urgent question about his way in the world. He already knows the answer, but it fails to satisfy him. So at great inconvenience, hardship, and even peril, he consults a powerful and refractory spirit who tries to evade his question by turning into anything in the world. Then, when the spirit sees he cannot get free of the man, and only then, he answers the man's question, not simply with a commonplace but with the same commonplace the man had been dissatisfied with before. Satisfied or not, however, the man now obeys the advice given him.

A foolish story? All the same, it is to be observed that Menelaus did get home. And it was a heroic thing to have hung onto Proteus through those terrifying changes and compelled him to be himself and answer up. Nor does it matter in the least to the story that Menelaus personally may have been a disagreeable old fool as well as a cuckold.

A poet also has one great and simple question, simple though it may take many forms indeed. Geoffrey Chaucer put it as well as anyone could, and in three lines at that:

What is this world? what asketh men to have?
Now with his love, now in his colde grave,
Allone, with-outen any companye.

(*"The Knight's Tale"*)

And a poet gets the simple answer he might expect, the one the world grudgingly gives to anyone who asks such a question: The world is this way, not that way, and you ask for more than you will be given, which the poet, being scarcely more fool than his fellowmen, knew already. But on the path from question to answer, hanging onto the slippery dissembler and shape-shifter Proteus, he will see many marvels; he will follow the metamorphoses of things in the metamorphoses of their phrases, and he will be so elated and ecstatic in this realm of wonders that the voice in which he speaks these things, down even to the stupid, obvious, and commonplace answer, will be to his hearers a solace and a happiness in the midst of sorrows:

When I do count the clock that tells the time,
And see the brave day sunk in hideous night;
When I behold the violet past prime,
And sable curls, all silver'd o'er with white;
When lofty trees I see barren of leaves,
Which erst from heat did canopy the herd,
And summer's green all girded up in sheaves,
Borne on the bier with white and bristly beard,
Then of thy beauty do I question make,
That thou among the wastes of time must go,
Since sweets and beauties must themselves forsake
And die as fast as they see others grow;

And nothing 'gainst Time's scythe can make defence
Save breed, to brave him when he takes thee hence.

(Shakespeare, Sonnet 12.)

Like Menelaus, the poet asks a simple question, to which, moreover, he already knows the unsatisfying answer. Question and answer, one might say, have to be present, although of themselves they seem to do nothing much; but they assert the limits of a journey to be taken. They are the necessary but not sufficient conditions of what really seems to matter here, the Protean encounter itself, the grasping and hanging on to the powerful and refractory spirit in its slippery transformations of a single force flowing through clock, day, violet, graying hair, trees dropping their leaves, the harvest in which, by a peculiarly ceremonial transmutation, the grain man lives by is seen without contradiction as the corpse he comes to. As for the answer to the question, it is not surprising nor meant to be surprising; it is only just.

On this point—that the answer comes as no surprise—poets show an agreement that quite transcends the differences of periods and schools. Alexander Pope's formula, "What oft was thought, but ne'er so well express'd," sometimes considered as the epitome of a shallow and parochial decorum, is not in essence other than this offered by John Keats:

I think Poetry should surprise by a fine excess, and not by Singularity—it should strike the Reader as a wording of his own highest thoughts, and appear almost a Remembrance. (Letter to John Taylor, 1818.)

In the present century, Robert Frost is strikingly in agreement:

A word about recognition: In literature it is our business to give people the thing that will make them say, "Oh yes I know what you mean." It is never to tell them something they don't know, but something they know and hadn't thought of saying. It must be something they recognize. (Letter to John Bartlett, in *Modern Poetics*, ed. James Scully, 1965.)

And the poet and critic John Crowe Ransom gives the thought a cryptically and characteristically elegant variation: "Poetry is the kind of knowledge by which we must know that we have arranged that we shall not know otherwise." Perhaps this point about recognition might be carried further, to the extreme at which it would be seen to pose the problem of how poetry, which at its highest has always carried, at least implicitly, a kind of Platonism and claimed to give, if not knowledge itself, what was more important, a "form" to knowledge, can survive the triumph of scientific materialism and a positivism minded to skepticism about everything in the world except its own self (where it turns credulous, extremely). The poet's adjustment, over two or three centuries, to a Newtonian cosmos, Kantian criticism, the spectral universe portrayed by physics has conspicuously not been a happy one and has led alternately or simultaneously to the extremes of rejection of reason and speaking in tongues on the one hand and the hysterical claim that poetry will save the world on the other. But of this let the Protean parable speak as it will.

There is another part to the story of Menelaus and Proteus, for Menelaus asked another question: What happened to my friends who were with me at Troy? Proteus replies, "Son of Atreus, why enquire too closely of me on this? To know or learn what I know about it is not your need: I warn you that when you hear all the truth your tears will not be far behind . . ." But he tells him all the same: "Of those others many went under; many came through . . ." And Menelaus does indeed respond with tears of despair, until Proteus advises him to stop crying and get started on the journey home. So it sometimes happens in poetry, too: the sorrowful contemplation of what is, consoles, in the end, and heals, but only after the contemplative process has been gone through and articulated in the detail of its change:

When to the sessions of sweet silent thought
I summon up remembrance of things past,
I sigh the lack of many a thing I sought,
And with old woes new wail my dear time's waste;
Then can I drown an eye, unused to flow,
For precious friends hid in death's dateless night,
And weep afresh love's long since cancell'd woe,
And moan the expense of many a vanish'd sight,
Then can I grieve at grievances foregone,

The poetic
question
and answer

And heavily from woe to woe tell o'er
 The sad account of fore-bemoaned moan,
 Which I new pay as if not paid before.
 But if the while I think on thee, dear friend,
 All losses are restor'd and sorrows end.
 (Shakespeare, Sonnet 30)

The nature
 of poetry

This poem, acknowledged to be a masterpiece by so many generations of readers, may stand as an epitome and emblem for the art altogether, about which it raises a question that must be put, although it cannot be satisfactorily and unequivocally answered: the question of whether poetry is a sacrament or a confidence game or both or neither. To reply firmly that poetry is not religion and must not promise what religion does is to preserve a useful distinction; nevertheless, the religions of the world, if they have nothing else in common, seem to be based on collections of sacred poems. Nor, at the other extreme, can any guarantee that poetry is not a confidence game be found in the often-heard appeal to the poet's "sincerity." One will never know whether Shakespeare wept all over the page while writing the 30th sonnet, though one inclines to doubt it, nor would it be to his credit if he did, nor to the reader's that he should know it or care to know it.

For one thing, the sonnet is obviously artful—that is, full of artifice—and even the artifice degenerates here and there into being artsy. "Then can I drown an eye, unused to flow." Surely that is poesy itself, at or near its worst, where the literal and the conventional, whatever their relations may have been for Shakespeare and the first reader of these sugar'd sonnets among his friends, now live very uncomfortably together (Ben Jonson's "Drink to me only with thine eyes" is a like example of this bathetic crossing of levels), though perhaps it has merely become unattractive as a result of changing fashions in diction.

Moreover, while the whole poem is uniquely Shakespearean, the bits and pieces are many of them common property of the age, what one writer called "joint stock company poetry." And the tricks are terribly visible, too; art is not being used to conceal art in such goings-on as "grieve at grievances" and "fore-bemoaned moan." "He who thus grieves will excite no sympathy," as Samuel Johnson sternly wrote of John Milton's style in the elegy "Lycidas," "he who thus praises will confer no honour."

Nor is that the worst of it. This man who so powerfully works on the reader's sympathies by lamenting what is past contrives to do so by thinking obsessively about litigation and, of all things, money; his hand is ever at his wallet, bidding adieu. He cannot merely "think" sweet silent thoughts about the past; no, he has to turn them into a court in "session," whereto he "summons" the probable culprit "remembrance"; when he "grieves," it is at a "grievance"—in the hands of the law again; finally, as with the sinners in Dante's *Divine Comedy*, his avarice and prodigality occupy two halves of the one circle: he bemoans his expenses while paying double the asking price.

And still, for all that, the poem remains beautiful; it continues to move both the young who come to it still innocent of their dear time's waste and the old who have sorrows to match its sorrows. As between confidence game and sacrament there may be no need to decide, as well as no possibility of deciding: elements of play and artifice, elements of true feeling, elements of convention both in the writing and in one's response to it, all combine to veil the answer. But the poem remains.

Effect of
 the poem
 as a whole

If it could be plainly demonstrated by the partisans either of unaided reason or revealed religion that poetry was metaphorical, mythological, and a delusion, while science, say, or religion or politics were real and true, then one might throw poetry away and live honestly though poorly on what was left. But, for better or worse, that is not the condition of man's life in the world; and perhaps men care for poetry so much—if they care at all—because, at last, it is the only one of man's many mythologies to be aware, and to make him aware, that it, and the others, are indeed mythological. The literary critic I.A. Richards, in a deep and searching consideration of this matter, concludes: "It is the privilege of poetry to preserve us from mistaking our notions either for things or for ourselves. Poetry is the completest mode of utterance."

The last thing Proteus says to Menelaus is strange indeed: You are not to die in Argos of the fair horse-pastures, not there to encounter death: rather will the Deathless Ones carry you to the Elysian plain, the place beyond the world. . . . There you will have Helen to yourself and will be deemed of the household of Zeus.

So the greatest of our poets have said, or not so much said, perhaps, as indicated by their fables, though nowadays people mostly sing a different tune. To be as the gods, to be rejoined with the beloved, the world forgotten. . . . Sacrament or con game? Homer, of course, is only telling an old story and promises mankind nothing; that is left to the priests to do; and in that respect poetry, as one critic puts it, must always be "a ship that is wrecked on entering the harbor." And yet the greatest poetry sings always, at the end, of transcendence; while seeing clearly and saying plainly the wickedness and terror and beauty of the world, it is at the same time humming to itself, so that one overhears rather than hears: All will be well. (H.Ne.)

Prosody

As it has come to be defined in modern criticism, the term prosody encompasses the study of all of the elements of language that contribute toward acoustic and rhythmic effects, chiefly in poetry but also in prose. The term derived from an ancient Greek word that originally meant a song accompanied by music or the particular tone or accent given to an individual syllable. Greek and Latin literary critics generally regarded prosody as part of grammar: it concerned itself with the rules determining the length or shortness of a syllable, with syllabic quantity, and with how the various combinations of short and long syllables formed the metres (*i.e.*, the rhythmic patterns) of Greek and Latin poetry. Prosody was the study of metre and its uses in lyric, epic, and dramatic verse. In sophisticated modern criticism, however, the scope of prosodic study has been expanded until it now concerns itself with what the 20th-century poet Ezra Pound called "the articulation of the total sound of a poem."

Prose as well as verse reveals the use of rhythm and sound effects; however, critics do not speak of "the prosody of prose" but of prose rhythm. The English critic George Saintsbury wrote *A History of English Prosody from the Twelfth Century to the Present* (3 vol., 1906–10), which treats English poetry from its origins to the end of the 19th century; but he dealt with prose rhythm in an entirely separate work, *A History of English Prose Rhythm* (1912). Many prosodic elements such as the rhythmic repetition of consonants (alliteration) or of vowel sounds (assonance) occur in prose; the repetition of syntactical and grammatical patterns also generates rhythmic effect. Traditional rhetoric, the study of how words work, dealt with acoustic and rhythmic techniques in Classical oratory and literary prose. But although prosody and rhetoric intersected, rhetoric dealt more exactly with verbal meaning than with verbal surface. Rhetoric dealt with grammatical and syntactical manipulations and with figures of speech; it categorized the kinds of metaphor. Modern critics, especially those who practice the New Criticism, might be considered rhetoricians in their detailed concern with such devices as irony, paradox, and ambiguity. These subjects are discussed at greater length below in the section *Literary criticism* and in the article RHETORIC.

Prosody
 and prose
 rhythm

This section considers prosody chiefly in terms of the English language—the only language that all of the readers of this article may be assumed to know. Some examples are given in other languages to illustrate particular points about the development of prosody in those languages; because these examples are pertinent only for their rhythm and sound, and not at all for their meaning, no translations are given. A further general discussion of the development of prosodic elements will be found in the section above, *Poetry*.

ELEMENTS OF PROSODY

As a part of modern literary criticism, prosody is concerned with the study of rhythm and sound effects as they occur in verse and with the various descriptive, historical,

and theoretical approaches to the study of these structures.

Scansion. The various elements of prosody may be examined in the aesthetic structure of prose. The celebrated opening passage of Charles Dickens' novel *Bleak House* (1853) affords a compelling example of prose made vivid through the devices of rhythm and sound:

Fog everywhere. Fog up the river, where it flows among green aits and meadows; fog down the river, where it rolls defiled among the tiers of shipping, and the waterside pollutions of a great (and dirty) city. Fog on the Essex Marshes, fog on the Kentish heights. Fog creeping into the cabooses of collier-brigs; fog lying out on the yards, and hovering in the rigging of great ships; fog drooping on the gunwales of barges and small boats. Fog in the eyes and throats of ancient Greenwich pensioners, wheezing by the firesides of their wards; fog in the stem and bowl of the afternoon pipe of the wrathful skipper . . .

Two phrases of five syllables each ("Fog everywhere"; "Fog up the river") establish a powerful rhythmic expectation that is clinched in repetition:

. . . fog down the river . . . Fog on the Essex . . . fog on the Kentish . . . Fog creeping into . . . fog drooping on the . . .

This phrase pattern can be scanned; that is, its structure of stressed and unstressed syllables might be translated into visual symbols:

' u u u u
 Fog down the ri ver.

(This scansion notation uses the following symbols: the acute accent ['] to mark metrically stressed syllables; the breve [˘] to mark metrically weak syllables; a single line [|] to mark the divisions between feet [*i.e.*, basic combinations of stressed and unstressed syllables]; a double line [||] to mark the caesura, or pause in the line; a rest [^] to mark a syllable metrically expected but not actually occurring.) Such a grouping constitutes a rhythmic constant, or cadence, a pattern binding together the separate sentences and sentence fragments into a long surge of feeling. At one point in the passage, the rhythm sharpens into metre; a pattern of stressed and unstressed syllables falls into a regular sequence:

' u u u u u u u u u u u u
 Fog on the | Es sex | mar shes, || fog on the |
 Ken tish | heights.

The line is a hexameter (*i.e.*, it comprises six feet), and each foot is either a dactyl (˘˘˘) or a trochee (˘˘).

The passage from Dickens is strongly characterized by alliteration, the repetition of stressed consonantal sounds:

Fog creeping into the cabooses of collier-brigs;

and by assonance, the patterned repetition of vowel sounds:

. . . fog down the river, where it rolls defiled among . . .

Here the vowel sounds are symmetrically distributed: short, long and long, short. Thus, it is clear that Dickens uses loosely structured rhythms, or cadences, an occasional lapse into metre, and both alliteration and assonance.

The rhythm and sound of all prose are subject to analysis; but compared with even the simplest verse, the "prosodic" structure of prose seems haphazard, unconsidered. The poet organizes his structures of sound and rhythm into rhyme, stanzaic form, and, most importantly, metre. Indeed, the largest part of prosodical study is concerned with the varieties of metre, the nature and function of rhyme, and the ways in which lines of verse fall into regular patterns or stanzas. An analysis of "Vertue" by the 17th-century English poet George Herbert reveals how the elements of prosody combine into a complex organism, a life sustained by the technical means available to the poet. When the metre is scanned with the symbols, it can be seen (and heard) how metre in this poem consists of the regular recurrence of feet, how each foot is a pattern of phonetically stressed and unstressed syllables.

' u u u u u u u u u u u u
 1 Sweet day, | so cool, | so calm, | so bright,
 ' u u u u u u u u u u u u
 2 The bri | dall of | the earth | and skie:
 ' u u u u u u u u u u u u
 3 The dew | shall weep | thy fall | to-night:

u u u u u u u u u u u u
 4 For thou | must die.

u u u u u u u u u u u u
 5 Sweet rose, | whose hue | an grie | and brave

u u u u u u u u u u u u
 6 Bids the | rash ga | zer wipe | his eye:

u u u u u u u u u u u u
 7 Thy root | is ev | er in | its grave,

u u u u u u u u u u u u
 8 And thou | must die.

u u u u u u u u u u u u
 9 Sweet spring, | full of | sweet dayes | and ro | ses.

u u u u u u u u u u u u
 10 A box | where sweets | com pac | ted lie:

u u u u u u u u u u u u
 11 My mu | sick shows | ye have | your clo | ses.

u u u u u u u u u u u u
 12 And all | must die.

u u u u u u u u u u u u
 13 Onely | a sweet | and ver | tuous soul,

u u u u u u u u u u u u
 14 Like sea | son'd tim | ber, ne | ver gives:

u u u u u u u u u u u u
 15 But though | the whole | world turn | to coal,

u u u u u u u u u u u u
 16 Then chief | ly lives.

The basic prosodic units are the foot, the line, and the stanza. The recurrence of similar feet in a line determines the metre; here there are three lines consisting of four iambic feet (*i.e.*, of four units in which the common pattern is the iamb—an unstressed syllable followed by a stressed syllable), which are followed by a line consisting of two iambic feet. Thus the stanza or recurring set of lines consists of three iambic tetrameters followed by one iambic dimeter. The stanzaic form is clinched by the use of rhyme; in "Vertue" the first and third and second and fourth lines end with the same sequence of vowels and consonants: bright/night, skie/pie, brave/grave, eye/pie, etc. It should be observed that the iambic pattern (˘˘) is not invariable; the third foot of line 5, the first foot of line 6, the second foot of line 9, and the first foot of line 13 are reversals of the iambic foot or trochees (˘˘). These reversals are called substitutions; they provide tension between metrical pattern and meaning, as they do in these celebrated examples from Shakespeare:

u u u u u u u u u u u u
 To be, | or not | to be, || that is | the ques tion . . .

Hamlet

u u u u u u u u u u u u
 His sil | ver skin | laced with | his gol | den blood . . .

Macbeth

Meaning, pace, and sound. Scansion reveals the basic metrical pattern of the poem; it does not, however, tell everything about its prosody. The metre combines with other elements, notably propositional sense or meaning, pace or tempo, and such sound effects as alliteration, assonance, and rhyme. In the fifth line of "Vertue," the reversed third foot occurring at "angry" brings that word into particular prominence; the disturbance of the metre combines with semantic reinforcement to generate a powerful surge of feeling. Thus, the metre here is expressive. The pace of the lines is controlled by the length of number of syllables and feet, line 5 obviously takes longer to read or recite. The line contains more long vowel sounds:

Sweet rose, whose hue angrie and brave . . .

Vowel length is called quantity. In English verse, quantity cannot by itself form metre although a number of English poets have experimented with quantitative verse. Generally speaking, quantity is a rhythmical but not a metrical feature of English poetry; it can be felt but it cannot be precisely determined. The vowel sounds in "Sweet rose" may be lengthened or shortened at will. No such options are available, however, with the stress patterns of words:

Quantity

Cadence
and metre

Scansion

the word an-gry cannot be read an-gry.

Assonance takes into account the length and distribution of vowel sounds. A variety of vowel sounds can be noted in this line:

Sweet day, so cool, so calm, so bright . . .

To borrow a term from music, the line modulates from *cē*, through *ā*, *ōō*, *ā*, to *ī*. Alliteration takes into account the recurrence and distribution of consonants:

so cool, so calm . . .

Sweet spring . . .

Rhyme normally occurs at the ends of lines; "Vertue" reveals, however, a notable example of interior rhyme, or rhyme within the line:

My musick shows ye have your closes . . .

Types of metre. *Syllable-stress metres.* It has been shown that the metre of "Vertue" is determined by a pattern of stressed and unstressed syllables arranged into feet and that a precise number of feet determines the measure of the line. Such verse is called syllable-stress verse (in some terminologies accentual-syllabic) and was the norm for English poetry from the beginning of the 16th century to the end of the 19th century. A line of syllable-stress verse is made up of either two-syllable (disyllabic) or three-syllable (trisyllabic) feet. The disyllabic feet are the iamb and the trochee (noted in the scansion of "Vertue"); the trisyllabic feet are the dactyl (˘˘˘) and anapest (˘˘˘).

Following are illustrations of the four principal feet found in English verse:

iambic	be hold
trochaic	ti ger
dactylic	des per ate
anapestic	un der stand

Some theorists also admit the spondaic foot (˘˘) and pyrrhic foot (˘˘) into their scansions; however, spondees and pyrrhics occur only as substitutions for other feet, never as determinants of a metrical pattern:

When to | the ses | sions of | sweet si | lent thought . . .

It has been noted that four feet make up a line of tetrameter verse; a line consisting of one foot is called monometer, of two dimeter, of three trimeter, of five pentameter, of six hexameter, and of seven heptameter. Lines containing more than seven feet rarely occur in English poetry.

The following examples illustrate the principal varieties of syllable-stress metres and their scansions:

Iambic (pentameter)

Then say | not Man's | im per | fect, || Heaven | in fault;
 Say ra | ther, || Man's | as per | fect as | he ought:
 His know | ledge meas | ured || to | his state | and place,
 His time | a mo | ment, || and | a point | his space.

Alexander Pope, *An Essay on Man* (1733-34)

Trochaic (dimeter)

Could I | catch that
 Nim ble | trai tor
 Scorn ful | Lau ra,
 Swift-foot | Lau ra,
 Soon then | would I
 Seek a | venge ment.

Thomas Campion (1602)

Dactylic (tetrameter)

After the | pangs of a | des per ate | Lover, ^
 When day and | night I have | sigh'd all in | vain, ^ ^
 Ah what a | pleas ure it | is to dis | co ver ^
 In her eyes | pi ty, who | cau ses my | pain! ^ ^

John Dryden, *An Evening's Love* (1671)

Anapestic (tetrameter)

The As syr | ian came down | like a wolf | on the fold,
 And his co | horts were gleam | ing in pur | ple and gold;
 And the sheen | of their spears | was like stars | on the sea,
 When the blue | wave rolls night | ly on deep | Ga li lee.

Lord Byron, "The Destruction of Sennacherib" (1815)

Syllable stress became more or less established in the poetry of Geoffrey Chaucer (c. 1340-1400). In the century that intervened between Chaucer and the early Tudor poets, syllable-stress metres were either ignored or misconstrued. By the end of the 16th century, however, the now-familiar iambic, trochaic, dactylic, and anapestic metres became the traditional prosody for English verse.

Strong-stress metres. In the middle of the 19th century, with Walt Whitman's free verse and Gerard Manley Hopkins' extensive metrical innovations, the traditional prosody was challenged. Antecedent to the syllable-stress metres was the strong-stress metre of Old English and Middle English poetry. Strong-stress verse is measured by count of stresses alone; the strong stresses are usually constant, but the number of unstressed syllables may vary considerably.

Strong-stress verse survives in nursery rhymes and children's counting songs:

One, two, || buckle my shoe;
 Three, four, || knock at the door;
 Five, six, || pick up sticks . . .

The systematic employment of strong-stress metre can be observed in the Old English epic poem *Beowulf* (c. 1000) and in William Langland's vision-poem, *Piers Plowman* ('A' Text, c. 1362):

In a somer sesun || whon softe was the sonne,
 I schop me in-to a schroud || a sheep as I were;
 In habite of an hermite || un-holy of werkes,

Wende I wydene in this world || wondres to here.

These lines illustrate the structural pattern of strong-stress metre. Each line divides sharply at the caesura (||), or medial pause; on each side of the caesura are two stressed syllables strongly marked by alliteration.

Strong-stress verse is indigenous to the Germanic languages with their wide-ranging levels of stressed syllables and opportunities for alliteration. Strong-stress metre was normative to Old English and Old Germanic heroic poetry, as well as to Old English lyric poetry. With the rising influence of French literature in the 12th and 13th centuries, rhyme replaced alliteration and stanzaic forms replaced the four-stress lines. But the strong-stress rhythm persisted; it can be felt in the anonymous love lyrics of the 14th century and in the popular ballads of the 15th century.

"Lord Randal" can be comfortably scanned to show a line of mixed iambic and anapestic feet; it clearly reveals, however, a four-stress structure:

'O where ha' you been, || Lord Randal, my son?

Challenges
to
traditional
prosody
in
the 19th
century

And where ha' you been, || my handsome young man?

'I ha' been at the greenwood; || mother, mak my bed soon.

For I'm wearied wi' huntin', || and fain wad lie down.'

A number of 20th-century poets, including Ezra Pound, T.S. Eliot, and W.H. Auden, have revived strong-stress metre. The versification of Pound's *Cantos* and Eliot's *Four Quartets* (1943) shows the vitality of the strong-stress, or, as they are often called, "native," metres.

Syllabic metres. Most of English poetry is carried by the strong-stress and syllable-stress metres. Two other kinds of metres must be mentioned: the purely syllabic metres and the quantitative metres. The count of syllables determines the metres of French, Italian, and Spanish verse. In French poetry the alexandrine, or 12-syllabled line, is a dominant metrical form:

O toi, qui vois la honte où je suis descendue,
 Implacable Vénus, suis-je assez confondue?
 Tu ne saurais plus loin pousser ta cruauté.
 Ton triomphe est parfait; tous tes traits ont porté.
Racine, Phèdre (1677)

Stress and pause in these lines are variable; only the count of syllables is fixed. English poets have experimented with syllabic metres; the Tudor poet Thomas Wyatt's translations from Petrarch's Italian poems of the 14th century attempted to establish a metrical form based on a decasyllabic or 10-syllabled line:

The long love that in my thought doth harbor,
 And in my heart doth keep his residence,
 Into my face presseth with bold pretense
 And there encampeth, spreading his banner.
 "The Lover for Shamefastness Hideth ..." (1557)

Most ears can detect that these lines waver between syllabic and syllable-stress metre: the second line falls into a pattern of iambic feet. Most ears also discover that the count of syllables alone does not produce any pronounced rhythmic interest; syllabic metres in English generate a prosody more interesting to the eye than to the ear.

Quantitative metres. Quantitative metres determine the prosody of Greek and Latin verse. Renaissance theorists and critics initiated a confused and complicated argument that tried to explain European poetry by the rules of Classical prosody and to draft laws of quantity by which European verse might move in the hexameters of the ancient Roman poets Virgil or Horace. Confusion was compounded because both poets and theorists used the traditional terminology of Greek and Latin prosody to describe the elements of the already existing syllable-stress metres; iambic, trochaic, dactylic, and anapestic originally named the strictly quantitative feet of Greek and Latin poetry. Poets themselves adapted the metres and stanzas of Classical poetry to their own languages; whereas it is not possible here to trace the history of Classical metres in European poetry, it is instructive to analyze some attempts to make English and German syllables move to Greek and Latin music. Because neither English nor German has fixed rules of quantity, the poets were forced to revise the formal schemes of the Classical paradigms in accordance with the phonetic structure of their own language.

A metrical paradigm much used by both Greek and Latin poets was the so-called Sapphic stanza. It consisted of three quantitative lines that scanned

— 0 — — 0 0 — 0 0 — 0
 — 0 0 — —

followed by a shorter line, called an Adonic,

"Sapphics" by the 19th-century English poet Algernon Charles Swinburne shows the Sapphic metre and stanza in English:

All the night sleep came not upon my eyelids,
 Shed not dew, nor shook nor unclosed a feather,
 Yet with lips shut close and with eyes of iron
 Stood and beheld me . . .
 Saw the white implacable Aphrodite,
 Saw the hair unbound and the feet unsandalled
 Shine as fire of sunset on western waters;
 Saw the reluctant . . .

The same metre and stanza in German are found in "Sapphische Ode," by the 19th-century poet Hans Schmidt, which was beautifully set to music by Johannes Brahms (Opus 94, No. 4):

Rosen brach ich nachts mir am dunklen Hage;
 süsser hauchten Duft sie, als je am Tage;
 doch verstreuten reich die bewegten Äste
 Tau den mich nässete.
 Auch der Küsse Duft mich wie nie berückte,
 die ich nachts vom Strauch deiner Lippen pflückte:
 doch auch dir, bewegt im Gemüt gleich jenem,
 tauten die Tränen.

Quantitative metres originated in Greek, a language in which the parts of speech appear in a variety of inflected forms (*i.e.*, changes of form to indicate distinctions in case, tense, mood, number, voice, and others). Complicated metrical patterns and long, slow-paced lines developed because the language was hospitable to polysyllabic metrical feet and to the alternation of the longer vowels characterizing the root syllables and the shorter vowels characterizing the inflected case-endings. The Classical metres can be more successfully adapted to German than to English because English lost most of its inflected forms in the 15th century, while German is still a highly inflected language. Thus Swinburne's "Sapphics" does not move as gracefully, as "naturally" as Schmidt's. A number of German poets, notably Goethe and Friedrich Hölderlin, both of the early 19th century, made highly successful use of the Classical metres. English poets, however, have never been able to make English syllables move in the ancient metres with any degree of comfort or with any sense of vital rhythmic force.

The American poet Henry Wadsworth Longfellow adapted the Classical hexameter for his *Evangeline* (1847):

This is the | for est pri | me val. The | mur mu ring | pines and
 the | hem locks . . .

In Virgil's *Aeneid*, Longfellow's Classical model, the opening line scans:

— 0 0 — 0 0 — — — 0 — 0 0 — —
 Ar ma vi rum que ca no, Troj ae qui pri mus ab or is

The rules determining length of syllable in Classical Greek and Latin poetry are numerous and complicated; they were established by precise grammatical and phonetic conventions. No such rules and conventions obtain in English; Robert Bridges, the British poet laureate and an authority on prosody, remarked in his *Poetical Works* (1912) that the difficulty of adapting English syllables to the Greek rules is "very great, and even deterrent." Longfellow's hexameter is in reality a syllable-stress line of five dactyls and a final trochee; syllabic quantity plays no part in determining the metre.

PROSODIC STYLE

The analysis of prosodic style begins with recognizing the metrical form the poet uses. Is he writing syllable-stress, strong-stress, syllabic, or quantitative metre? Or is he using a nonmetrical prosody? Again, some theorists would not allow that poetry can be written without metre; the examples of Whitman and many 20th-century innovators, however, have convinced most modern critics that a nonmetrical prosody is not a contradiction in terms but an obvious feature of modern poetry. Metre has not disappeared as an important element of prosody; indeed, some of the greatest poets of the modern period—William Butler Yeats, T.S. Eliot, Ezra Pound, Wallace Stevens—revealed themselves as masters of the traditional metres. They also experimented with newer prosodies based on prose cadences, on expansions of the blank-verse line, and revivals of old forms—such as strong-stress and ballad metres. Also noteworthy are the "visual" prosodies fostered by the poets of the Imagist movement and by such experimenters as e.e. cummings. Cummings revived the practice of certain 17th-century poets (notably George Herbert) of "shaping" the poem by typographic arrangements.

The prosodic practice of poets has varied enormously with the historical period, the poetic genre, and the poet's

The dominant alexandrine

Sapphics in English and German

Non-metrical prosody

individual style. In English poetry, for example, during the Old English period (to 1100), the strong-stress metres carried both lyric and narrative verse. In the Middle English period (from c. 1100 to c. 1500), stanzaic forms developed for both lyric and narrative verse. The influence of French syllable counting pushed the older stress lines into newer rhythms; Chaucer developed for *The Canterbury Tales* a line of 10 syllables with alternating accent and regular end rhyme—an ancestor of the heroic couplet. The period of the English Renaissance (from c. 1500 to 1660) marks the fixing of syllable-stress metre as normative for English poetry. Iambic metre carried three major prosodic forms: the sonnet, the rhyming couplet, and blank verse. The sonnet was the most important of the fixed stanzaic forms. The iambic pentameter rhyming couplet (later known as the heroic couplet) was used by Christopher Marlowe for his narrative poem *Hero and Leander* (1598); by John Donne in the early 17th century for his satires, his elegies, and his longer meditative poems. Blank verse (unrhymed iambic pentameter), first introduced into English in a translation by Henry Howard, earl of Surrey, published in 1557, became the metrical norm for Elizabethan drama. The period of the Renaissance also saw the refinement of a host of lyric and song forms; the rapid development of English music during the second half of the 16th century had a salutary effect on the expressive capabilities of poetic rhythms.

The personal element. A poet's choice of a prosody obviously depends on what his language and tradition afford; these are primary considerations. The anonymous author of the Old English poem *Deor* used the conventional four-stress metric available to him; but he punctuated groups of lines with a refrain:

paes ofereode: þisses swa maeg!
(that passed away: this also may!)

The refrain adds something to the prosodic conventions of regulated stress, alliteration, and medial pause: a sense of a smaller and sharper rhythmic unit within the larger rhythms of the given metre. While the poet accepts from history his language and from poetic convention the structure of his metre, he shapes his own style through individual modifications of the carrying rhythms. When critics speak of a poet's "voice," his personal tone, they are also speaking of his prosodic style.

Prosodic style must be achieved through a sense of tension; it is no accident that the great masters of poetic rhythm work against the discipline of a given metrical form. In his sonnets, Shakespeare may proceed in solemn iambic regularity, creating an effect of measured progression through time and its legacy of suffering and despair:

No longer mourn for me when I am dead
Than you shall hear the surly sullen bell
Give warning to the world that I am fled . . .
"Sonnet 71"

Or he may wrench the metre and allow the reader to feel the sudden violence of his feelings, the power of a conviction raised to a command:

Let me | not to | the mar | riage of | true minds
Ad mit | im pe | di ments. | Love is | not love . . .
"Sonnet 116"

The first two feet of the first line are trochaic reversals; the last two feet comprise a characteristic pyrrhic-spondaic formation. A trochaic substitution is quite normal in the first foot of an iambic pentameter line; a trochaic substitution in the second foot, however, creates a marked disturbance in the rhythm. There is only one "normal" iambic foot in the first line; this line runs over (or is enjambed) to the second line with its three consecutive iambic feet followed by a strong caesura and reversed fourth foot. These lines are, in Gerard Manley Hopkins' terms, metrically "counter-pointed"; trochees, spondees, and pyrrhics are heard against a ground rhythm of regular iambs. Without the ground rhythm, Shakespeare's expressive departures would not be possible.

A poet's prosodic style may show all of the earmarks of revolt against prevailing metrical practice. Whitman's

celebrated "free verse" marks a dramatic break with the syllable-stress tradition; he normally does not count syllables, stresses, or feet in his long sweeping lines. Much of his prosody is rhetorical: that is, Whitman urges his language into rhythm by such means as anaphora (*i.e.*, repetition at the beginning of successive verses) and the repetition of syntactical units. He derives many of his techniques from the example of biblical verses, with their line of various types of parallelism. But he often moves toward traditional rhythms; lines fall into conventional parameters:

O past! O happy life! O songs of joy!
"Out of the Cradle Endlessly Rocking" (1859)

Or they fall more often into dissyllabic hexameters:

Borne through the smoke of the battles and pierc'd
with missiles I saw them . . .
"When Lilacs Last in the Dooryard Bloom'd" (1865–66)

Despite the frequent appearance of regular metrical sequences, Whitman's lines cannot be scanned by the usual graphic method of marking syllables and feet; his prosody, however, is fully available to analysis. The shape on the page of the lines below (they comprise a single strophe or verse unit) should be noted, specifically the gradual elongation and sudden diminution of line length. Equally noteworthy are the repetition of the key word *carols*, the alliteration of the *s* sounds, and the use of words in falling (trochaic) rhythm, "lagging," "yellow," "waning":

Shake out carols!
Solitary here, the night's carols!
Carols of lonesome love! death's carols!
Carols under that lagging, yellow, waning moon!
O under that moon where she droops almost down into
the sea!
O reckless despairing carols.
"Out of the Cradle"

No regular metre moves these lines; but a clearly articulated rhythm—produced by shape, thematic repetitions, sound effects, and patterns of stress and pause—defines a prosody.

Whitman's prosody marks a clear break with previous metrical practices. Often a new prosody modifies an existing metrical form or revives an obsolete one. In "Gerontion" (1920), T.S. Eliot adjusted the blank-verse line to the emotionally charged, prophetic utterance of his persona, a spiritually arid old man:

After such knowledge, what forgiveness? Think now
History has many cunning passages, contrived corridors
And issues, deceives with whispering ambitions,
Guides us by vanities. Think now . . .
(From T.S. Eliot, *Collected Poems 1909–1962*,
Harcourt Brace Jovanovich, Inc.)

The first three lines expand the pentameter line beyond its normal complement of stressed and unstressed syllables; the fourth line contracts, intensifying the arc of feeling. Both Pound and Eliot used stress prosodies. Pound counted out four strong beats and used alliteration in his brilliant adaptation of the old English poem "The Seafarer" (1912):

Chill its chains are; chafing sighs
Hew my heart round and hunger begot
Mere-weary mood. Lest man known not
That he on dry land loveliest liveth . . .
(From Ezra Pound, *Personae*, Copyright 1926 by
Ezra Pound. Reprinted by permission of New Directions Publishing Corporation.)

He uses a similar metric for the energetic opening of his "Canto I." Eliot mutes the obvious elements of the form in the celebrated opening of *The Waste Land* (1922):

April is the cruellest month, breeding
Lilacs out of the dead land, mixing
Memory and desire, stirring
Dull roots with spring rain.
(From T.S. Eliot, *Collected Poems 1909–1962*,
Harcourt Brace Jovanovich, Inc.)

Here is the "native metre" with its falling rhythm, elegiac tone, strong pauses, and variably placed stresses. If this is free verse, its freedoms are most carefully controlled. "No verse is free," said Eliot, "for the man who wants to do a good job."

The prosodic styles of Whitman, Pound, and Eliot—

Whitman's
innova-
tions

Style
through
tension

The stress
prosodies
of Pound
and Eliot

The
prosody of
Tennyson
and
Browning

though clearly linked to various historical antecedents—are innovative expressions of their individual talents. In a sense, the prosody of every poet of genius is unique; rhythm is perhaps the most personal element of the poet's expressive equipment. Alfred Lord Tennyson and Robert Browning, English poets who shared the intellectual and spiritual concerns of the Victorian age, are miles apart in their prosodies. Both used blank verse for their dramatic lyrics, poems that purport to render the accents of real men speaking. The blank verse of Tennyson's "Ulysses" (1842) offers smoothly modulated vowel music, carefully spaced spondaic substitutions, and unambiguous pentameter regularity:

The long day wanes; the slow moon climbs; the deep
Moans round with many voices. Come, my friends,
'Tis not too late to seek a newer world.

Browning's blank verse aims at colloquial vigour; its "irregularity" is a function not of any gross metrical violation—it always obeys the letter of the metrical law—but of the adjustment of abstract metrical pattern to the rhythms of dramatic speech. If Tennyson's ultimate model is Milton's Baroque prosody with its oratorical rhythms, Browning's model was the quick and nervous blank verse of the later Elizabethan dramatists. Characteristic of Browning's blank verse are the strong accents, involuted syntax, pregnant caesuras, and headlong energy in "The Bishop Orders His Tomb at St. Praxed's Church" (1845):

Vanity, saith the preacher, vanity!
Draw round my bed: is Anselm keeping back?
Nephews—sons mine . . . ah God, I know not! Well—
She, men would have to be your mother once,
Old Gandolf envied me, so fair she was!

Influence of period and genre. In the lyric genres, the rhythms of the individual poet—or, in the words of the 20th-century American poet Robert Lowell, "the person himself"—can be heard in the prosody. In the long poem, the dramatic, narrative, and didactic genres, a period style is more likely to be heard in prosody. The blank-verse tragedy of the Elizabethan and Jacobean dramatists, the blank verse of Milton's *Paradise Lost* (1667) and its imitators in the 18th century (James Thomson and William Cowper), and the heroic couplet of Neoclassical satiric and didactic verse, each, in different ways, defines the age in which these prosodies flourished. The flexibility and energy of the dramatic verse of Marlowe, Shakespeare, and John Webster reflect the later Renaissance with its nervous open-mindedness, its obsessions with power and domination, and its lapses into despair. Miltonic blank verse, based on Latin syntax and adaptations of the rules of Latin prosody, moved away from the looseness of the Elizabethans and Jacobean toward a more ceremonial style. It is a Baroque style in that it exploits the musical qualities of sounds for their ornamental values. The heroic couplet, dominating the poetry of the entire 18th century, was unequivocally a prosodic period style; its elegance and epigrammatic precision entirely suited an age that valued critical judgment, satiric wit, and the powers of rationality.

Prosody in
dramatic
verse

It is in dramatic verse, perhaps, that a prosody shows its greatest vitality and clarity. Dramatic verse must make a direct impression not on an individual reader able to reconsider and meditate on what he has read but on an audience that must immediately respond to a declaiming actor or a singing chorus. The ancient Greek dramatists developed two distinct kinds of metres: "stichic" forms (*i.e.*, consisting of "stichs," or lines, as metrical units) such as the iambic trimeter for the spoken dialogues; and lyric, or strophic, forms (*i.e.*, consisting of stanzas), of great metrical intricacy, for the singing and chanting of choruses. Certain of the Greek metres developed a particular ethos; characters of low social standing never were assigned metres of the lyric variety. Similar distinctions obtained in Elizabethan drama. Shakespeare's kings and noblemen speak blank verse; comic characters, servants, and country bumpkins discourse in prose; clowns, romantic heroines, and supernatural creatures sing songs. In the early tragedy *Romeo and Juliet*, the chorus speaks in "excellent conceited" sonnets: in what was one of the most popular and easily recognized lyric forms of the period.

The metrical forms used by ancient and Renaissance dramatists were determined by principles of decorum. The use or non-use of a metrical form (or the use of prose) was a matter of propriety; it was important that the metre be suitable to the social status and ethos of the individual character as well as be suitable to the emotional intensity of the particular situation. Decorum, in turn, was a function of the dominant Classical and Neoclassical theories of imitation.

THEORIES OF PROSODY

Ancient critics like Aristotle and Horace insisted that certain metres were natural to the specific poetic genres; thus, Aristotle (in the *Poetics*) noted that "Nature herself, as we have said, teaches the choice of the proper measure." In epic verse the poet should use the heroic measure (dactylic hexameter) because this metre most effectively represents or imitates such qualities as grandeur, dignity, and high passion. Horace narrowed the theory of metrical decorum, making the choice of metre prescriptive; only an ill-bred and ignorant poet would treat comic material in metres appropriate to tragedy. Horace prepared the way for the legalisms of the Renaissance theorists who were quite willing to inform practicing poets that they used "feete without joyntes," in the words of Roger Ascham, Queen Elizabeth's tutor, and should use the quantitative metres of Classical prosody.

The Middle Ages. During the Middle Ages little of importance was added to actual prosodic theory: in poetic practice, however, crucial developments were to have important ramifications for later theorists. From about the second half of the 6th century to the end of the 8th century, Latin verse was written that no longer observed the rules of quantity but was clearly structured on accentual and syllabic bases. This change was aided by the invention of the musical sequence; it became necessary to fit a musical phrase to a fixed number of syllables, and the older, highly complex system of quantitative prosody could not be adapted to simple melodies that must be sung in sequential patterns. In the musical sequence lies the origin of the modern lyric form.

The
musical
sequence

The 9th-century hymn "Ave maris stella" is a striking instance of the change from quantitative to accentual syllabic prosody: each line contains three trochaic feet determined not by length of syllable but by syllabic intensity or stress:

Ave maris stella
Dei mater alma
atque semper virgo,
felix caeli porta.
Sumens illud Ave
Gabrielis ore,
funda nos in pace,
mutans nomen Evae.

The rules of quantity have been disregarded or forgotten; rhyme and stanza and a strongly felt stress rhythm have taken their place. In the subsequent emergence of the European vernacular literatures, poetic forms follow the example of the later Latin hymns. The earliest art lyrics, those of the Provençal troubadours of the 12th and 13th centuries, show the most intricate and ingenious stanzaic forms. Similarly, the Goliardic songs of the *Carmina Burana* (13th century) reveal a rich variety of prosodic techniques; this "Spring-song" embodies varying lines of trochees and iambs and an ababedced rhyme scheme:

Ver redit optatum
Cum gaudio,
Flore decoratum
Purpureo;
Aves edunt cantus
Quam dulciter!
Revirescit nemus,
Cantus est amoenus
Totaliter.

The Renaissance. Renaissance prosodic theory had to face the fact of an accomplished poetry in the vernacular that was not written in metres determined by "rules" handed down from the practice of Homer and Virgil. Nevertheless, the classicizing theorists of the 16th century made a determined attempt to explain existing poetry

Rules of
"trewe
versifying"

by the rules of short and long and to draft "laws" by which modern verse might move in Classical metres. Roger Ascham, in *The Scholemaster* (1570), attacked "the Gothic . . . barbarous and rude Ryming" of the early Tudor poets. He admitted that Henry Howard, earl of Surrey, did passably well as a poet but complained that Surrey did not understand "perfitte and trewe versifying"; that is, Surrey did not compose his English verses according to the principles of Latin and Greek quantitative prosody.

Ascham instigated a lengthy argument, continued by succeeding theorists and poets, on the nature of English prosody. Sir Philip Sidney, Gabriel Harvey, Edmund Spenser, and Thomas Campion all (to use Saintsbury's phrase) committed whoredom with the enchantment of quantitative metric. While this hanky-panky had no adverse effect on poetry itself (English poets went on writing verses in syllable-stress, the prosody most suitable to the language), it produced misbegotten twins of confusion and discord, whose heirs, however named, are still apparent today. Thus, those who still talk about "long and short" (instead of stressed and unstressed), those who perpetuate a punitive prosodic legalism, and those who regard prosody as an account of what poets should have done and did not, trace their ancestry back to Elizabethan dalliance and illicit classicizing.

Although Renaissance prosodic theory produced scarcely anything of value to either literary criticism or poetic technique—indeed, it did not even develop a rational scheme for scanning existing poetry—it raised a number of important questions. What were the structural principles animating the metres of English verse? What were the aesthetic nature of prosody and the functions of metre? What were the connections between poetry and music? Was poetry an art of imitation (as Aristotle and all of the Neoclassical theorists had maintained), and was its sister art painting; or was poetry (as Romantic theory maintained) an art of expression, and prosody the element that produced (in Coleridge's words) the sense of musical delight originating (in T.S. Eliot's words) in the auditory imagination?

Pope's
doctrine of
imitation

The 18th century. Early in the 18th century, Pope affirmed, in his *Essay on Criticism* (1711), the classic doctrine of imitation. Prosody was to be more nearly onomatopoeic; the movement of sound and metre should represent the actions they carry:

'Tis not enough no harshness gives offence,
The sound must seem an Echo to the sense:
Soft is the strain when Zephyr gently blows,
And the smooth stream in smoother numbers flows;
But when loud surges lash the sounding shoar,
The hoarse, rough verse should like the torrent roar.
When Ajax strives some rock's vast weight to throw,
The line too labours, and the words move slow;
Not so, when swift Camilla scours the plain,
Flies o'er th' unbending corn, and skims along the main.

In 18th-century theory the doctrine of imitation was joined to numerous strictures on "smoothness," or metrical regularity. Theorists advocated a rigid regularity; minor poets composed in a strictly regular syllable-stress verse devoid of expressive variations. This regularity itself expressed the rationalism of the period. The prevailing dogmas on regularity made it impossible for Samuel Johnson to hear the beauties of Milton's versification; he characterized the metrically subtle lines of "Lycidas" as "harsh" and without concern for "numbers." Certain crosscurrents of metrical opinion in the 18th century, however, moved toward new theoretical stances. Joshua Steele's *Prosodia Rationalis* (1779) is an early attempt to scan English verse by means of musical notation. (A later attempt was made by the American poet Sidney Lanier in his *Science of English Verse*, 1880.) Steele's method is highly personal, depending on an idiosyncratic assigning of such musical qualities as pitch and duration to syllabic values; but he recognized that a prosodic theory must take into account not merely metre but "all properties or accidents belonging to language." His work foreshadows the current concerns of the structural linguists who attempt an analysis of the entire range of acoustic elements contributing to prosodic effect. Steele is also the first "timer" among metrists; that is, he

bases his scansion on musical pulse and claims that English verse moves in either common or triple time. Modern critics of musical scanners have pointed out that musical scansion constitutes a performance, not an analysis of the metre, that it allows arbitrary readings, and that it levels out distinctions between poets and schools of poetry.

The 19th century. With the Romantic movement and its revolutionary shift in literary sensibility, prosodic theory became deeply influenced by early 19th-century speculation on the nature of imagination, on poetry as expression—"the spontaneous overflow of powerful feelings," in Wordsworth's famous phrase—and on the concept of the poem as organic form. The discussion between Wordsworth and Coleridge on the nature and function of metre illuminates the crucial transition from Neoclassical to modern theories. Wordsworth (in his "Preface" to the *Lyrical Ballads*, 1800) followed 18th-century theory and saw metre as "superadded" to poetry; its function is more nearly ornamental, a grace of style and not an essential quality. Coleridge saw metre as being organic; it functions together with all of the other parts of a poem and is not merely an echo to the sense or an artifice of style. Coleridge also examined the psychologic effects of metre, the way it sets up patterns of expectation that are either fulfilled or disappointed:

As far as metre acts in and for itself, it tends to increase the vivacity and susceptibility both of the general feelings and of the attention. This effect it produces by the continued excitement of surprize, and by the quick reciprocations of curiosity still gratified and still re-excited, which are too slight indeed to be at any one moment objects of distinct consciousness, yet become considerable in their aggregate influence. As a medicated atmosphere, or as wine during animated conversation; they act powerfully, though themselves unnoticed. Where, therefore, correspondent food and appropriate matter are not provided for the attention and feelings thus roused, there must needs be a disappointment felt; like that of leaping in the dark from the last step of a staircase, when we had prepared our muscles for a leap of three or four.

Biographia Literaria, XVIII (1817)

Romantic literary theory, although vastly influential in poetic practice, had little to say about actual metrical structure. Coleridge described the subtle relationships between metre and meaning and the effects of metre on the reader's unconscious mind; he devoted little attention to metrical analysis. Two developments in 19th-century poetic techniques, however, had greater impact than any prosodic theory formulated during the period. Walt Whitman's nonmetrical prosody and Gerard Manley Hopkins' far-ranging metrical experiments mounted an assault on the traditional syllable-stress metric. Both Whitman and Hopkins were at first bitterly denounced, but, as is often the case, the heresies of a previous age become the orthodoxies of the next. Hopkins' "sprung rhythm"—a rhythm imitating natural speech, using mixed types of feet and counterpointed verse—emerged as viable techniques in the poetry of Dylan Thomas and W.H. Auden. It is virtually impossible to assess Whitman's influence on the various prosodies of modern poetry. Such American poets as Hart Crane, William Carlos Williams, and Theodore Roethke all have used Whitman's long line, extended rhythms, and "shaped" strophes.

The 20th century. Since 1900 the study of prosody has emerged as an important and respectable part of literary study. George Saintsbury published his great *History of English Prosody* during the years 1906–10. Sometime later, a number of linguists and aestheticians turned their attention to prosodic structure and the nature of poetic rhythm. Graphic prosody (the traditional syllable and foot scansion of syllable-stress metre) was placed on a securer theoretical footing. A number of prosodists, taking their lead from the work of Joshua Steele and Sidney Lanier, have recently attempted to use musical notation to scan English verse. For the convenience of synoptic discussion, modern prosodic theorists may be divided into four groups: the linguists who examine verse rhythm as a function of phonetic structures; the aestheticians who examine the psychologic effects, the formal properties, and the phenomenology of rhythm; the musical scanners, or "timers," who try to adapt the procedures of musical notation to

The poem
as organic
form

Linguists,
aestheti-
cians,
timers, and
tradition-
alists

metrical analysis; and the traditionalists who rely on the graphic description of syllable and stress to uncover metrical paradigms. It is necessary to point out that only the traditionalists concern themselves specifically with metrical form; aestheticians, linguists, and timers all examine prosody in its larger dimensions.

Modern structural linguistics has placed the study of language on a solid scientific basis. Linguists have measured the varied intensities of syllabic stress and pitch and the durations of junctures or the pauses between syllables. These techniques of objective measurement have been applied to prosodic study. The Danish philologist Otto Jespersen's early essay "Notes on Metre" (1900) made a number of significant discoveries. He established the principles of English metre on a demonstrably accurate structural basis: he recognized metre as a gestalt phenomenon (*i.e.*, with emphasis on the configurational whole); he saw metrics as descriptive science rather than proscriptive regulation. Jespersen's essay was written before the burgeoning interest in linguistics, but since World War II numerous attempts have been made to formulate a descriptive science of metrics.

It has been noted that Coleridge defined metrical form as a pattern of expectation, fulfillment, and surprise. Taking his cue from Coleridge, the British aesthetician I.A. Richards in *Principles of Literary Criticism* (1924) developed a closely reasoned theory of the mind's response to rhythm and metre. His theory is organic and contextual; the sound effects of prosody have little psychologic effect by themselves. It is prosody in conjunction with "its contemporaneous other effects"—chiefly meaning or propositional sense—that produces its characteristic impact on our neural structures. Richards insists that everything that happens in a poem depends on the organic environment: in his *Practical Criticism* (1929) he constructed a celebrated "metrical dummy" to "support [an] argument against anyone who affirms that the mere sound of verse has *independently* any considerable aesthetic virtue." For Richards the most important function of metre is to provide aesthetic framing and control; metre makes possible, by its stimulation and release of tensions, "the most difficult and delicate utterances."

Other critics, following the Neo-Kantian theories of the philosophers Ernst Cassirer and Susanne Langer, have suggested that rhythmic structure is a species of symbolic form. Harvey Gross in *Sound and Form in Modern Poetry* (1964) saw rhythmic structure as a symbolic form, signifying ways of experiencing organic processes and the phenomena of nature. The function of prosody, in his view, is to image life in a rich and complex way. Gross's theory is also expressive; prosody articulates the movement of feeling in a poem. The unproved assumption behind Gross's expressive and symbolic theory is that rhythm is in some way iconic to human feeling; that a particular rhythm or metre symbolizes, as a map locates the features of an actual terrain, a particular kind of feeling.

Rhythmic structure as a symbolic form

The most sophisticated argument for musical seanson is given by Northrop Frye in his influential *Anatomy of Criticism* (1957). He differentiates between verse that shows unmistakable musical quality and verse written according to the imitative doctrines current in the Renaissance and Neoclassic periods. All of the poetry written in the older strong-stress metric, or poetry showing its basic structure, is musical poetry, and its structure resembles the music contemporary with it.

The most convincing case for traditional "graphic prosody" has been made by the American critics W.K. Wimsatt and Monroe C. Beardsley. Their essay "The Concept of Meter" (1965) argues that both the linguists and musical scanners do not analyze the abstract metrical pattern of poems but only interpret an individual performance of the poem. Poetic metre is not generated by any combination of stresses and pauses capable of precise scientific measurement; rather, metre is generated by an abstract pattern of syllables standing in positions of relative stress to each other. In a line of iambic pentameter

Preserved in Milton's or in Shakespeare's name . . .

the "or" of the third foot is only slightly stronger than the preceding syllable "-ton's," but this very slight difference makes the line recognizable as iambic metre. Wimsatt and Beardsley underline the paradigmatic nature of metre; as an element in poetic structure, it is capable of exact abstraction.

Non-Western theories. The metres of the verse of ancient India were constructed on a quantitative basis. A system of long and short syllables, as in Greek, determined the variety of complicated metrical forms that are found in poetry of post-Vedic times—that is, after the 5th century BC.

Chinese prosody is based on the intricate tonal system of the language. In the T'ang dynasty (AD 618–907) the metrical system for classical verse was fixed. The various tones of the language were subsumed under two large groups, even tones and oblique tones. Patterned arrangements of tones and the use of pauses, or caesuras, along with rhyme determine the Chinese prosodic forms.

Japanese poetry is without rhyme or marked metrical structure; it is purely syllabic. The two main forms of syllabic verses are the tanka and the haiku. Tanka is written in a stanza of 31 syllables that are divided into alternating lines of five and seven syllables. Haiku is an extremely concentrated form of only 17 syllables. Longer poems of 40 to 50 lines are also written; however, alternate lines must contain either five or seven syllables. The haiku form has been adapted to English verse and has become in recent years a popular form. Other experimenters in English syllabic verse show the influence of Japanese prosody. Syllabic metre in English, however, is limited in its rhythmic effects; it is incapable of expressing the range of feeling that is available in the traditional stress and syllable-stress metres. (Ha.G.)

Japanese prosody

NARRATIVE FICTION

Epic

An ambiguous term, "epic" is used most often to designate a long narrative poem recounting heroic deeds, though it has also been loosely used to describe novels, such as Tolstoy's *War and Peace*, and motion pictures, such as Eisenstein's *Ivan the Terrible*. In literary usage, the term encompasses both oral and written compositions. The prime examples of the oral epic are Homer's *Iliad* and *Odyssey*. Outstanding examples of the written epic include Virgil's *Aeneid* and Lucan's *Pharsalia* in Latin; *Chanson de Roland* in medieval French; Ariosto's *Orlando Furioso* and Tasso's *Gerusalemme liberata* in Italian; *Poema* (or *Cantar*) *de mio Cid* in Spanish; and Milton's *Paradise Lost* and Spenser's *Faerie Queene* in English. There are also seriocomic epics, such as the *Morgante* of a 15th-century Italian poet, Pulci, and the pseudo-Homeric *Battle of the Frogs and Mice*. Another distinct group is made up of the

so-called beast epics—narrative poems written in Latin in the Middle Ages and dealing with the struggle between a cunning fox and a cruel and stupid wolf. Underlying all of the written forms is some trace of an oral character, partly because of the monumental persuasiveness of Homer's example but more largely because the epic was, in fact, born of an oral tradition. It is on the oral tradition of the epic form that this article will focus.

GENERAL CHARACTERISTICS

An epic may deal with such various subjects as myths, heroic legends, histories, edifying religious tales, animal stories, or philosophical or moral theories. Epic poetry has been used by peoples all over the world and in different ages to transmit their traditions from one generation to another, without the aid of writing. These traditions frequently consist of legendary narratives about the glorious deeds of their national heroes. Thus scholars have often

The epic as a genre of heroic ages

identified "epic" with a certain kind of heroic oral poetry, which comes into existence in so-called heroic ages. Such ages have been experienced by many nations, usually at a stage of development in which they have had to struggle for a national identity. This effort, combined with such other conditions as an adequate material culture and a sufficiently productive economy, tend to produce a society dominated by a powerful and warlike nobility, constantly occupied with martial activities, whose individual members seek, above all, everlasting fame for themselves and for their lineages.

Uses of the epic. The main function of poetry in heroic-age society appears to be to stir the spirit of the warriors to heroic actions by praising their exploits and those of their illustrious ancestors, by assuring a long and glorious recollection of their fame, and by supplying them with models of ideal heroic behaviour. One of the favorite pastimes of the nobility in heroic ages in different times and places has been to gather in banquet halls to hear heroic songs, in praise of famous deeds sung by professional singers as well as by the warriors themselves. Heroic songs also were often sung before a battle, and such recitations had tremendous effect on the morale of the combatants. Among the Fulani (Fulbe) people in The Sudan, for instance, whose epic poetry has been recorded, a nobleman customarily set out in quest of adventures accompanied by a singer (*mabo*), who also served as his shield bearer. The singer was thus the witness of the heroic deeds of his lord, which he celebrated in an epic poem called *baudi*.

The aristocratic warriors of the heroic ages were thus members of an illustrious family, a link in a long chain of glorious heroes. And the chain could snap if the warrior failed to preserve the honour of the family, whereas, by earning fame through his own heroism, he could give it new lustre. Epic traditions were to a large extent the traditions of the aristocratic families: the Old French word *geste*, used for a form of epic that flourished in the Middle Ages, means not only a story of famous deeds but also a genealogy.

The passing of a heroic age does not necessarily mean the end of its heroic oral poetry. An oral epic tradition usually continues for as long as the nation remains largely illiterate. Usually it is after the heroic age has passed that the narratives about its legendary heroes are fully elaborated. Even when the nobility that originally created the heroic epic perishes or loses interest, the old songs can persist as entertainments among the people. Court singers, then, are replaced by popular singers, who recite at public gatherings. This popular tradition, however, must be distinguished from a tradition that still forms an integral part of the culture of a nobility. For when a heroic epic loses its contact with the banquet halls of the princes and noblemen, it cannot preserve for long its power of renewal. Soon it enters what has been called the reproductive stage in the life cycle of an oral tradition, in which the bards become noncreative reproducers of songs learned from older singers. Popular oral singers, like the *guslari* of the Balkans, no doubt vary their songs to a certain extent each time they recite them, but they do so mainly by transposing language and minor episodes from one acquired song to another. Such variations must not be confounded with the real enrichment of the tradition by succeeding generations of genuine oral poets of the creative stage. The spread of literacy, which has a disastrous effect on the oral singer, brings about a quick corruption of the tradition. At this degenerate stage, the oral epic soon dies out if it is not written down or recorded.

The ancient Greek epic exemplifies the cycle of an oral tradition. Originating in the late Mycenaean period, the Greek epic outlasted the downfall of the typically heroic-age culture (c. 1100 BC) and maintained itself through the "Dark Age" to reach a climax in the Homeric poems by the close of the Geometric period (900–750 BC). After Homer, the activity of the *aidoi*, who sang their own epic songs at the courts of the nobility, slowly declined. During the first half of the 7th century, the *aidoi* produced such new poems as those of Hesiod and some of the earlier poems of what was to become known as the Epic Cycle. Between 625 and 575 BC the *aidoi* gave way to

oral reciters of a new type, called rhapsodes or "stitchers of songs," who declaimed for large audiences the already famous works of Homer while holding in their hand a staff (*rhabdos*), which they used to give emphasis to their words. It seems probable that these rhapsodes, who played a crucial role in the transmission of the Homeric epic, were using some sort of written aids to memory before Homeric recitations were adopted in 6th-century Athens as part of the Panathenaic festivals held each year in honour of the goddess Athena.

Verbal formulas. To compose and to memorize long narrative poems like the *Iliad* and the *Odyssey*, oral poets used a highly elaborate technical language with a large store of traditional verbal formulas, which could describe recurring ideas and situations in ways that suited the requirements of metre. So long as an oral epic tradition remains in its creative period, its language will be continually refined by each generation of poets in opposite directions, refinements that are called scope and economy. Scope is the addition of new phrases to express a larger number of recurrent concepts in varying metrical values fitting the possible positions in a verse. Economy is the elimination of redundancies that arise as gifted poets invent new set phrases that duplicate, both in a general sense and in metrical value, the formulas that already exist in the traditional stock.

Nowhere has this refinement proceeded any nearer to perfection than in the language of the Homeric epic. As has been shown by statistical analysis, it exhibits a remarkable efficiency, both in the rareness of unnecessarily duplicative variants and in the coverage of each common concept by the metrical alternatives useful in the composition of the six-foot metric line the Greeks used for epic poetry.

Thus, for example, if the idea of a ship has to be expressed at the end of a line of verse, the ship may be described as "well-trimmed" (*nēos eisēs*), "curved" (*nēos amphielissa*), or "dark-prowed" (*nēos kyanoproiros*), depending entirely on the number of feet that remain to be filled by the phrase in the hexameter; if the phrase has to cover the two final feet of the verse and the words have to be put in the dative case, the formula "of a well-trimmed ship" will be replaced by "to a black ship" (*nēi melainē*). The sole occurrence of "Zeus who gathers lightning" (*steropēgereta Zeus*), which is an exact metrical equivalent of the more common "Zeus who delights in thunder" (*Zeus terpikeraunos*), constitutes one of the very few actual duplications of such formulas found in Homer.

Finally, some of the typical scenes in the heroic life, such as the preparation of a meal or sacrifice or the launching or beaching of a ship, contain set descriptions comprising several lines that are used by rote each time the events are narrated.

This highly formalized language was elaborated by generations of oral poets to minimize the conscious effort needed to compose new poems and memorize existing ones. Because of it, an exceptionally gifted *aidos*, working just prior to the corruption of the genre, could orally create long and finely structured poems like the *Iliad* and the *Odyssey*, and those poems could then be transmitted accurately by the following generations of rhapsodes until complete written texts were produced.

BASES

Oral heroic poetry, at its origin, usually deals with outstanding deeds of kings and warriors who lived in the heroic age of the nation. Since the primary function of this poetry is to educate rather than to record, however, the personages are necessarily transformed into ideal heroes and their acts into ideal heroic deeds that conform to mythological or ideological patterns. Some of these patterns are archetypes found all over the world, while others are peculiar to a specific nation or culture. Thus, in many epic traditions, heroes are born as a result of an illegitimate union of a maiden mother with a divine or supernatural being; they are exposed at birth, fed by an animal, and brought up by humble foster parents in a rustic milieu; they grow up with marvellous speed, fight a dragon—in their first combat—to rescue a maiden whom

The popular tradition of the epic

Heroic patterns and deeds

they marry, and die young in circumstances as fabulous as those that surrounded their birth.

In the traditions of Indo-European peoples a hero is often a twin, who acquires soon after his supernatural birth an invulnerability that has one defect, generally of his heel or of some other part of his foot, which ultimately causes his death. He is educated by a blacksmith, disguises himself as a woman at some time in his youth, and conquers a three-headed dragon, or some other kind of triple opponent, in his first battle. He then begets, by a foreign or supernatural woman, a child who, reared by his mother in her country, becomes a warrior as brave as his father. When this child meets his unknown father, the latter fails to recognize him, so that the father kills his own child after a long and fierce single combat. The hero, himself, usually dies after committing the third of three sins.

In Japan, to take another example, renowned members of the warrior aristocracy of the past, who have acquired the status of popular heroes, are in many cases supplied in their legend with four exceptionally brave and faithful retainers called their *shi-temō*, the guardians of the four cardinal points; these form the closest entourage of their lord—who is usually depicted as excelling in command but not in physical strength—and defend him from dangers. The retainers reflect a mythological model, taken from Buddhism, of four *deva* kings, who guard the teaching of the Buddha against the attack of the devils.

A striking pattern for a number of epic traditions has been found in a so-called "tripartite ideology" or "trifunctional system" of the Indo-Europeans. The concept was based on the discovery of the remarkable philosophy of a prehistoric nation that survived as a system of thought in the historic Indo-European civilizations and even in the subconsciousness of the modern speakers of Indo-European tongues.

This philosophy sees in the universe three basic principles that are realized by three categories of people: priests, warriors, and producers of riches. In conformity with this philosophy, most Indo-European epics have as their central themes interaction among these three principles or functions which are: (1) religion and kingship; (2) physical strength; (3) fecundity, health, riches, beauty, and so forth. In the long Indian epic the *Mahābhārata*, for example, the central figures, the Pāṇḍava brothers, together with their father Paṇḍu, their two uncles Dhṛtarāṣṭra and Vidura, and their common wife, Draupadi, correspond to traditional deities presiding over the three functions of the Indo-European ideology.

During the first part of their earthly career, the Pāṇḍava suffer constantly from the persistent enmity and jealousy of their cousins, Duryodhana and his 99 brothers, who, in reality, are incarnations of the demons Kāli and the Paulastya. The demons at first succeeded in snatching the Kingdom from the Pāṇḍava and in exiling them. The conflict ends in a devastating war, in which all the renowned heroes of the time take part. The Pāṇḍava survive the massacre, and establish on earth a peaceful and prosperous reign, in which Dhṛtarāṣṭra and Vidura also participate.

This whole story, it has been shown, is a transposition to the heroic level of an Indo-European myth about the incessant struggle between the gods and the demons since the beginning of the world. Eventually, it results in a bloody eschatological battle, in which the gods and the devils exterminate each other. The destruction of the former world order, however, prepares for a new and better world, exempt from evil influences, over which reign a few divine survivors of the catastrophe.

EARLY PATTERNS OF DEVELOPMENT

In the ancient Middle East. The earliest-known epic poetry is that of the Sumerians. Its origin has been traced to a preliterate Heroic Age, not later than 3000 bc, when the Sumerians had to fight, under the direction of a warlike aristocracy, for possession of this fertile Mesopotamian land. Among the extant literature of this highly gifted people are fragments of narrative poems recounting the heroic deeds of their early kings: Enmerkar, Lugalbanda, and Gilgamesh. By far the most important in the development of Mesopotamian literature are the five poems

of the Gilgamesh cycle. This epic tells the odyssey of a king, Gilgamesh, part human and part divine, who seeks immortality. A god who dislikes his rule, fashions a wild man, Enkidu, to challenge him. Enkidu first lives among wild animals, then goes to the capital and engages in a trial of strength with Gilgamesh, who emerges victorious. The two, now friends, set out on various adventures, in one of which they kill a wild bull that the goddess of love had sent to destroy Gilgamesh because he spurned her marriage proposal. Enkidu dreams the gods have decided he must die for the death of the bull, and, upon awakening, he does fall ill and die. Gilgamesh searches for a survivor of the Babylonian flood to learn how to escape death. The survivor shows him where to find a plant that renews youth, but after Gilgamesh gets the plant it is snatched away by a serpent. Gilgamesh returns, saddened, to his capital.

The legend of Gilgamesh was taken over by the Babylonians, who developed it into a long and beautiful poem, one of the masterpieces of mankind.

Another Babylonian epic, composed around 2000 bc, is called in Akkadian *Enuma elish*, after its opening words, meaning "When on high." Its subject is not heroic but mythological. It recounts events from the beginning of the world to the establishment of the power of Marduk, the great god of Babylon. The outline of a Babylonian poem narrating the adventure of a hero named Adapa ("Man") can be reconstructed from four fragmentary accounts. It shares with the *Epic of Gilgamesh* the theme of man's loss of an opportunity for immortality.

Among clay tablets of the 14th century bc, covered with inscriptions in an old Phoenician cuneiform alphabet, from Ras Shamra (the site of ancient Ugarit), in northern Syria, there are important fragments of three narrative poems. One of these is mythological and recounts the career of the god Baal, which seems to coincide with the yearly cycle of vegetation on earth. As was usual with the death of gods in the ancient Mediterranean world, Baal's end brings about a drought that ceases only with his resurrection. Another fragment, about a hero named Aqhat, is perhaps a transposition of this myth of Baal to the human level. Just as the death of Baal is avenged on his slayer by Baal's sister Anath, so is the murder of Aqhat, which also causes a drought, revenged by his sister Paghat. Since the end of the poem is missing, however, it is not known whether Paghat, like Anath, succeeded in bringing her brother back to life.

The third fragment, the Ugaritic epic of Keret, has been interpreted as a Phoenician version of the Indo-European theme of the siege of an enemy city for the recovery of an abducted woman. This theme is also the subject of the Greek legend of the Trojan War and of the Indian epic *Rāmāyaṇa*. The fragmentary text does not reveal, however, whether the expedition of Keret, like that of the Achaean army against Troy, was meant to regain the hero's wife or to acquire for him a new bride.

The Greek epic. In its originative stage, especially, the Greek epic may have been strongly influenced by these oriental traditions. The Greek world in the late Bronze Age was related to the Middle East by so many close ties that it formed an integral part of the Levant. At Ugarit a large quarter of the city was occupied by Greek merchants, whose presence is also attested, among other places, at the gate of Mesopotamia, at Alalakh, in what is now Turkey. Thus, it is no surprise that, for example, the Greek myth about the succession of the divine kingship told in the *Theogony* of Hesiod and elsewhere is paralleled in a Hittite version of a Hurrian myth. In it, Anu, Kumarbi, and the storm god respectively, parallel Uranus, Cronos, and Zeus in the *Theogony*. The Hittites had continuous diplomatic relations with the Achaeans of Greece, whose princes went to the royal court at Hattusa to perfect their skill with the chariot. The Greeks, therefore, had ample opportunity to become familiar with Hittite myths.

The *Epic of Gilgamesh* was then well-known in the Levant, as is indicated by discoveries of copies of it throughout this wide area. Many parallels with the *Epic of Gilgamesh* have been pointed out in the *Odyssey*: the encounters of Odysseus with Circe and Calypso on their

The legend of Gilgamesh

Central themes of Indo-European epics

Eastern influences on the Greek epic

mythical isles, for instance, closely resemble the visit by Gilgamesh to a divine woman named Siduri, who keeps an inn in a marvellous garden of the sun god near the shores of ocean. Like the two Greek goddesses, Siduri tries to dissuade Gilgamesh from the pursuit of his journey by representing the pleasures of life, but the firm resolution of the hero obliges her finally to help him cross the waters of death. In the *Iliad*, Patroclus, who dies as a substitute for his king and dearest friend, Achilles, and then gives Achilles a description of the miserable condition of man after his death, bears striking similarities to the friend of Gilgamesh, Enkidu.

If these are indeed borrowings, it is all the more remarkable that they are used in Homer to express a view of life and a heroic temper radically different from those of the Sumerian epic of Mesopotamia. Gilgamesh persists in his quest of immortality even when Siduri shows him the vanity of such an ambition, but Odysseus shuns a goddess's offer of everlasting life, preferring to bear his human condition to the end. The loss of a beloved friend does not make Achilles seek desperately to escape from death; instead he rushes into combat to revenge Patroclus, although he knows that he is condemning himself to an early death, and that the existence of a king in Hades will be incomparably less enviable than that of a slave on earth. The Mesopotamian mind never tires of expressing man's deep regret at not being immortal through stories about ancient heroes who, despite their superhuman strength and wisdom, and their intimacy with gods, failed to escape from death. A decisively different idea, however, is fundamental to the Greek heroic view of life. It has been demonstrated that the Greek view is derived from an Indo-European notion of justice—that each being has a fate (*moira*) assigned to him and marked clearly by boundaries that should never be crossed. Man's energy and courage should, accordingly, be spent not in exceeding the proper limits of his human condition but in bearing it with style, pride, and dignity, gaining as much fame as he can within the boundaries of his *moira*. If he is induced by Folly (Ate, personified as a goddess of mischief) to commit an excess (*hybris*) with regard to his *moira*, he will be punished without fail by the divine vengeance personified as Nemesis.

At the beginning of the *Iliad*, a plague decimates the Achaean army because its commander in chief, Agamemnon, refuses to return a captive, Chryseis, to her father, a priest of Apollo who offers a generous ransom. By unjustly insulting Achilles, Agamemnon commits another excess that causes the defeat of his army. Achilles, in the meantime, lets Ate take possession of his mind and refuses, to the point of excess, to resume his fight. He thus brings about a great misfortune, the loss of his dearest companion, Patroclus. Patroclus, however, also contributes to his own death by his *hybris* in pushing his triumph too far, ignoring Achilles' order to come back as soon as he has repulsed the enemy far from the Greek ships. The death of Hector also results from his *hybris* in rejecting the counsel of Polydamas and maintaining his army on the plain after the return of Achilles to combat. After so many disasters caused by the mischievous action of Ate among men, the last book of the *Iliad* presents a noble picture of Priam and Achilles, who submit piously to the orders of Zeus, enduring with admirable courage and moderation their respective fates.

On the other hand, at the beginning of the *Odyssey*, Zeus evokes the ruin that Aegisthus will have to suffer for having acted "beyond his due share" by marrying Clytemnestra and murdering Agamemnon. This sets an antithesis to the story of the wise Odysseus, who, to accomplish his destiny as a mortal hero, never changes his purpose trying always to make the best of his countless misfortunes. He earns by this the favour of Athena and succeeds eventually in regaining Ithaca and punishing the wooers of Penelope for their *hybris* during his long absence. Present scholarship inclines to the view that such admirably well-structured poems as the *Iliad* and the *Odyssey* could have been created only by a single highly gifted poet whose name was Homer. This position contrasts with the extreme skepticism that marked all phases

of Homeric criticism during the previous century. Yet the personality of Homer remains unknown and nothing certain is known about his life.

In comparison, information derived from his own works is fairly plentiful about the other great epic poet of Greece, Hesiod. He produced them presumably around 700 BC, while tilling a farm in Askra, a small village of Boeotia. The social and geographical background of his poems, called didactic because of their occasionally moral and instructive tone, differs from the aristocratic society of Ionian Asia Minor that Homer addressed. Despite their different style, subjects, and view of life, however, Hesiod's *Theogony* and the *Works and Days* illustrate the same basic conception of justice as the Homeric epic. The *Theogony* describes a long sequence of primordial events that resulted in the present world order, in which man's inescapable lot is assigned to him by Zeus. The *Works and Days* explains, through a series of three myths, why the lot of man is to work hard to produce riches. Man has to shut his ears to the goddess who causes wars and lawsuits, listening only to the goddess who urges him to toil more laboriously than his neighbour to become richer. Pain and suffering have become unavoidable since Pandora opened the fatal jar containing all the ills of mankind at Prometheus' house in conformity with the will of Zeus. Moreover, the age of the race of iron has arrived when the fate of human beings is not to pass their lives in perpetual banquets or warfare, as did the preceding races, but to suffer constantly the fatigue and misery of labour. As long as the goddesses Aidos (a personification of the sense of shame) and Nemesis (a personification of divine retribution) stay with mankind, however, helping people observe their *moira* without committing excesses, man can still gain riches, merits, and glory by the sweat of his brow. Only if he knows how to avoid all faults in doing his daily work will he not offend Justice (Dikē), the sensitive virgin daughter of Zeus. This is why it is so vitally important for a farmer to know all the rules listed in the rest of the poem about seemingly trivial details of his work.

LATER VARIATIONS

The Latin epic. Latin epic poetry was initiated in the 3rd century BC by Livius Andronicus, who translated the *Odyssey* into the traditional metre of Saturnian verse. It was not until the 1st century BC, however, that Rome possessed a truly national epic in the unfinished *Aeneid* of Virgil (70–19 BC), who used Homer as his model. The story of Aeneas' journey, recounted in the first six books, is patterned after the *Odyssey*, with many imitative passages and even direct translations, while the description of the war in the last six books abounds with incidents modelled after those from the *Iliad*. More basically, however, Virgil made use of another model, Rome's own national legend about the war fought under Romulus against the Sabines. This legend preserves, in a historical disguise, an original Indo-European myth about a primitive conflict between the gods of sovereignty and war and the gods of fecundity, ending with the unification of the two divine races. In the development of this theme by Virgil, Aeneas and the Etruscans can be seen as representing the gods of sovereignty and war, and the Latins representing the gods of fecundity. Aeneas, who has brought the Trojan gods to Rome, is forced to fight with the help of the Etruscans against the Latins. It is the destiny of Aeneas to rule, and it is the fate of the Latins to share their land and women with the invaders and to accept Aeneas as their king. This resembles the unification of the warring races that climaxes the Indo-European myth.

The power exercised by the Indo-European ideological pattern on the Roman mind even under the empire is seen in the *Pharsalia* of Lucan (AD 39–65). In this historical epic, Cato, Caesar, and Pompey are depicted respectively as moral, warlike, and popular in a way that gives the story a clear trifunctional structure.

Germanic epics. A typical Heroic Age occurred during the wanderings of the Germanic tribes from the 3rd to the 6th centuries AD. Out of this, too, came a rich oral tradition, from which developed in the Middle Ages many epic poems. One of the greatest of these is the Old English

Hesiod's
contribution

The Greek
heroic view
of life

The persistence of
the Indo-European
pattern

Beowulf, written down in the 8th century. Archetypal Indo-European themes also reappear in these epics. For example, the theme of the fatal fight between father and son is recounted in the German *Hildebrandslied*, of which a 67-line fragment is extant. Again, an heroic version of the Indo-European myth about the rescue of the Sun Maiden from her captivity by the Divine Twins, which also provided the basic plot of the Greek Trojan cycle and the Indian *Rāmāyaṇa*, is found in the German *Gudrun* (c. 1230).

Chansons de geste. The French chansons de geste are epic poems whose action takes place during the reign of Charlemagne and his immediate successors. The *Chanson de Roland*, probably written down about the end of the 11th century, is by far the most refined of the group. The story of the poem had developed from a historical event, the annihilation of the rear guard of Charlemagne's army at Roncesvalles in the Pyrenees in 778 by Basque mountaineers. The Basques, however, are transformed in the epic to the Saracens, who to a later generation typified France's enemies in Spain. The other chansons de geste, none of which is comparable to *Roland* as a literary work, have been classified into three main cycles. The cycle of Guillaume d'Orange forms a biography of William (probably a historical, count of William of Toulouse, who had, like the hero of the epic, a wife called Guibourg and a nephew, Vivien, and who became a monk in 806). Guibourg, the most faithful of wives, and the noble Vivien take prominent roles in the epic. The so-called Cycle of the Revolted Knights groups those poems that tell of revolts of feudal subjects against the emperor (Charlemagne or, more usually, his son, Louis). The Cycle of the King consists of the songs in which Charlemagne himself is a principal figure.

Arthurian Romance. The Arthurian Romance seems to have developed first in the British Isles, before being taken to the Continent by Bretons, who migrated to Brittany in the 6th and 7th centuries. The core of the legend about Arthur and his knights derives from lost Celtic mythology. Many of the incidents in the former parallel the deeds of such legendary Irish characters as Cú Chulainn, an Ulster warrior said to have been fathered by the god Lug, and Finn, hero of the Fenian cycle about a band of warriors defending Ireland, both of whom are gods transformed into human heroes. The earliest extant works on Arthurian themes are four poems of Chrétien de Troyes, written in French between 1155 and 1185: *Erec*, *Yvain*, *Le Chevalier de la Charette* (left unfinished by Chrétien and completed by Godefroy de Lagny), and an unfinished *Perceval*. In German, after 1188, Hartmann von Aue (who also wrote two legendary poems not belonging to the Arthurian cycle, *Gregorius* and *Poor Henry*) modelled his *Erec* and *Iwein* on those of Chrétien. The story of *Perceval* was given a full account by Wolfram von Eschenbach (c. 1170–1220) in his *Parzival* and in the unfinished *Titarel*. Another incomplete work of Wolfram, *Willehalm*, deals with the legend of William of Orange. *Tristan* of Gottfried von Strassburg is based directly on the older French version of Thomas of Britain (c. 1170–80). The romance proper, however, although it has similarities to the epic, differs in its lack of high purpose: fictions are told for their entertainment value rather than as models for national heroism. Developed in France in the Middle Ages, the romance is usually an adventure story with a strong love interest, intimately associated with the "courtly love" tradition of that time. (For further treatment, see below *Romance*.)

The epic in Japan. In Japan, there were in ancient times families of reciters (*katari-be*) whose duty was to hand down myths and legends by word of mouth and to narrate them during official ceremonies and banquets. After the introduction of Chinese letters, however, from the 4th century AD onward, these traditional tales were put in writing and the *katari-be* professional gradually died out. By the end of the 7th century, each clan of the ruling aristocracy seems to have possessed a written document that recounted the mythology and legendary history of Japan in a form biased in favour of the clan concerned. These family documents were collected at the command of the emperor Temmu (672–686) and were used as basic

materials for the compilation of the first national chronicles of Japan, the *Koji-ki* (712) and the *Nihon shoki* (720). The myths and legends that are contained in the earlier parts of these two books derive, therefore, from the oral tradition of *katari-be*. Although no document preserves those narrations in their primitive form, it is generally assumed that they were originally in the form of poems. Many scholars believe that they were genuine epic poems, which were produced during a period of incessant warfare around the 4th century. At that time mounted aristocratic warriors of the future imperial family struggled to extend its power over the larger part of Japan. Exploits of warriors, such as the emperor Jimmu or Prince Yamato-Takeru, in the earliest extant texts—the *Koji-ki* and *Nihon shoki* of the 8th century—probably derive from a heroic epic about the wars of conquest of the first emperors, whose legendary feats were transformed into those of a few idealized heroic figures.

The middle of the Heian period (794–1185) saw the emergence of a new class of warrior known as samurai. They attached a greater importance to fame than to life. The battles they fought became the subject of epic narratives that were recited by itinerant blind priests to the accompaniment of a lute-like instrument called a *biwa*.

In the early part of the 13th century, tales about the wars of the preceding century, fought between the two strongest families of samurai, the Genji, or Minamoto, and the Heike, or Taira, were compiled in three significant war chronicles. The *Hōgen monogatari* and the *Heiji monogatari* deal with two small wars, the Hōgen (1156) and Heiji (1159), in which the Genji and Heike warriors fought for opposing court factions. The structure of the two works is roughly the same. Each celebrates the extraordinary prowess of a young Genji warrior, Minamoto Tametomo in the *Hōgen monogatari* and Minamoto Yoshihira in the *Heiji monogatari*; each hero fights to the finish in exemplary manner not so much to win, for from the beginning each foresees the defeat of his own side, as for the sake of fame; and the consummate courage of the two heroes forms a striking contrast to the cowardice of court aristocrats. The bitterly fought Gempei War (1180–85), in which survivors of the Genji family challenged and defeated the Heike, is recounted in detail in the *Heike monogatari*, the greatest epic of Japanese literature. The sudden decline and ultimate extinction of the proud Heike, whose members had held the highest offices of the imperial court, illustrates the Buddhist philosophy of the transitory nature of all things: it invites the readers to seek deliverance from the world of sufferings through a faith that will take them to a land of eternal felicity at the moment of their death. The work is filled with tales of heroic actions of brave warriors. The most conspicuous is Minamoto Yoshitsune, one of the chief commanders of the Genji army; the legend of this man of military genius continued to develop in later literature, so that he has become the most popular hero of Japanese legend.

The later written epic. The vitality of the written epic is manifested by such masterworks as the Italian *Divine Comedy* of Dante (1265–1321) and the great Portuguese patriotic poem *Os Lusíadas* of Luiz de Camões (1524–80), which celebrates the voyage of Vasco da Gama to India. In more recent times, novels and long narrative poems written by such major authors as Scott, Byron, Tennyson, William Morris, and Melville were patterned, to some extent, on the epic. Their fidelity to the genre, however, is found primarily in their large scope and their roots in a national soil; their distance from the traditional oral epic tends to be considerable.

Among the epics written in modern times, the Finnish *Kalevala* (first ed. 1835; enlarged ed. 1849) occupies a very special position. This is because its author, the Finnish poet-scholar Elias Lönnrot (1802–84), who composed this masterpiece by combining short popular songs (*runot*) collected by himself among the Finns, had absorbed his material so well and identified himself so completely with the *runo* singers. He thus came close to showing what the oral epic, which he could study only at its degenerative stage, might have been at its creative stage, on the lips of an exceptionally gifted singer. (A.V.)

Distinctions between the epic and the romance

Tales of the samurai

Fable, parable, and allegory

Fables, parables, and allegories are forms of imaginative literature or spoken utterance constructed in such a way that their readers or listeners are encouraged to look for meanings hidden beneath the literal surface of the fiction. A story is told or perhaps enacted whose details—when interpreted—are found to correspond to the details of some other system of relations (its hidden, allegorical sense). The poet, for example, may describe the ascent of a hill in such a way that each physical step corresponds to a new stage in the soul's progress toward a higher level of existence.

Many forms of literature elicit this kind of searching interpretation, and the generic term for the cluster is allegory; under it may be grouped fables, parables, and other symbolic shapings. Allegory may involve either a creative or an interpretive process: either the act of building up the allegorical structure and giving "body" to the surface narrative or the act of breaking down this structure to see what themes or ideas run parallel to it.

NATURE AND OBJECTIVES

Allegory and myth. The fate of allegory, in all its many variations, is tied to the development of myth and mythology. Every culture embodies its basic assumptions in stories whose mythic structures reflect the society's prevailing attitudes toward life. If the attitudes are disengaged from the structure, then the allegorical meaning implicit in the structure is revealed. The systematic discipline of interpreting the real meaning of a text (called the hermeneutic process) plays a major role in the teaching and defense of sacred wisdom, since religions have traditionally preserved and handed down the old beliefs by telling exemplary stories; these sometimes appear to conflict with a system of morality that has in the meantime developed, and so their "correct" meaning can only be something other than the literal narration of events. Every culture puts pressure on its authors to assert its central beliefs, which are often reflected in literature without the author's necessarily being aware that he is an allegorist. Equally, determined critics may sometimes find allegorical meaning in texts with less than total justification—instances might include the Hebraic-Christian mystical interpretation of the Old Testament's Song of Solomon, an erotic marriage poem, or the frequent allegorizing of classical and modern literature in the light of Freud's psychoanalytic discoveries. Some awareness of the author's intention seems necessary in order to curb unduly fanciful commentary.

The allegorical mode. The range of allegorical literature is so wide that to consider allegory as a fixed literary genre is less useful than to regard it as a dimension, or mode, of controlled indirectness and double meaning (which, in fact, all literature possesses to some degree). Critics usually reserve the term allegory itself for works of considerable length, complexity, or unique shape. Thus, the following varied works might be called allegories: the biblical parable of the sower; *Everyman*, the medieval morality play; *The Pilgrim's Progress*, by John Bunyan; Jonathan Swift's *Gulliver's Travels*; *The Scarlet Letter*, by Nathaniel Hawthorne; William Wordsworth's "Ode: Intimations of Immortality"; Nikolay Gogol's *Dead Souls*; *The Picture of Dorian Gray*, by Oscar Wilde; and the plays *Six Characters in Search of an Author*, by Luigi Pirandello; *Waiting for Godot*, by Samuel Beckett; and *Who's Afraid of Virginia Woolf?*, by Edward Albee. No one genre can take in such modal range.

Fable. Fable and parable are short, simple forms of naïve allegory. The fable is usually a tale about animals who are personified and behave as though they were humans. The device of personification is also extended to trees, winds, streams, stones, and other natural objects. The earliest of these tales also included humans and gods as characters, but fable tends to concentrate on animating the inanimate. A feature that isolates fable from the ordinary folktale, which it resembles, is that a moral—a rule of behaviour—is woven into the story.

Parable. Like fable, the parable also tells a simple story. But, whereas fables tend to personify animal characters—

often giving the same impression as does an animated cartoon—the typical parable uses human agents. Parables generally show less interest in the storytelling and more in the analogy they draw between a particular instance of human behaviour (the true neighbourly kindness shown by the good Samaritan in the Bible story, for example) and human behaviour at large. Parable and fable have their roots in preliterate oral cultures, and both are means of handing down traditional folk wisdom. Their styles differ, however. Fables tend toward detailed, sharply observed social realism (which eventually leads to satire), while the simpler narrative surface of parables gives them a mysterious tone and makes them especially useful for teaching spiritual values.

Derivation of the terms. The original meanings of these critical terms themselves suggest the direction of their development. Fable (from the Latin *fabula*, "a telling") puts the emphasis on narrative (and in the medieval and Renaissance periods was often used when speaking of "the plot" of a narrative). Parable (from Greek *parabolē*, a "setting beside") suggests a juxtaposition that compares and contrasts this story with that idea. Allegory (from Greek *allos* and *agoreuein*, an "other-speaking") suggests a more expanded use of deceptive and oblique language. (In early Greek, though, the term allegory itself was not used. Instead, the idea of a hidden, underlying meaning is indicated by the word *hyponoia*—literally, "underthought"—and this term is used of the allegorical interpretation of the Greek poet Homer.)

Diverse objectives. *Fable.* Fables teach a general principle of conduct by presenting a specific example of behaviour. Thus, to define the moral that "People who rush into things without using judgment run into strange and unexpected dangers," Aesop—the traditional "father" of the fable form—told the following story:

There was a dog who was fond of eating eggs. Mistaking a shell-fish for an egg one day, he opened his mouth wide and swallowed it down in one gulp. The weight of it in his stomach caused him intense pain. "Serve me right," he said, "for thinking that anything round must be an egg."

By a slight change of emphasis, the fabulist could have been able to draw a moral about the dangerous effects of gluttony.

Because the moral is embodied in the plot of the fable, an explicit statement of the moral need not be given, though it usually is. Many of these moral tag lines have taken on the status of proverb because they so clearly express commonly held social attitudes.

The Aesopian fables emphasize the social interactions of human beings, and the morals they draw tend to embody advice on the best way to deal with the competitive realities of life. With some irony, fables view the world in terms of its power structures. One of the shortest Aesopian fables says: "A vixen sneered at a lioness because she never bore more than one cub. 'Only one,' she replied, 'but a lion.'" Foxes and wolves, which the poet Samuel Taylor Coleridge called "Everyman's metaphor" for cunning and cruelty, appear often as characters in fables chiefly because, in the human world, such predatory cunning and cruelty are able to get around restraints of justice and authority. The mere fact that fables unmask the "beast in me," as James Thurber, the 20th-century American humorist and fabulist, put it, suggests their satirical force. Subversive topical satire in tsarist and Soviet Russia is often called "Aesopism"; all comic strips that project a message (such as the Charles Schulz creation "Peanuts" and Walt Kelly's "Pogo") have affinities with Aesop's method.

Parable. Parables do not analyze social systems so much as they remind the listener of his beliefs. The moral and spiritual stress of the form falls upon memory rather than on the critical faculty. The audience hearing the parable is assumed to share a communal truth but perhaps to have set it aside or forgotten it. The rhetorical appeal of a parable is directed primarily toward an elite, in that a final core of its truth is known only to an inner circle, however simple its narrative may appear on the surface (a number of the parables that Christ used for teaching, for example, conveyed figuratively the meaning of the elusive concept Kingdom of Heaven).

The moral objective of fables

Allegory. Allegory, as the basic process of arousing in the reader or listener a response to levels of meaning, provides writers with the structure of fables, parables, and other related forms. By awakening the impulse to question appearances and by bringing order to mythological interpretation, allegory imparts cultural values. A measure of allegory is present in literature whenever it emphasizes thematic content, ideas rather than events. Generally, the allegorical mode flourishes under authoritarian conditions. Thus it found sustenance during the age of medieval Christendom, when Christian dogma sought universal sway over the mind of Western man. As such, allegory was a means of freedom under conditions of strong restraint. In general, realism, mimetic playfulness, and the resistance to authority tend to counteract the allegorical process, by loosening its stratified forms. This unbinding of symbolic hierarchies has forced allegory to seek new structures in the modern period. Nevertheless, through allegorical understanding, the great myths continue to be reread and reinterpreted, as the human significance of the new interpretations is passed down from one generation to the next. The abiding impression left by the allegorical mode is one of indirect, ambiguous, even enigmatic symbolism, which inevitably calls for interpretation.

Diversity of forms. Since an allegorical purpose can inform works of literature in a wide range of genres, it is not surprising to find that the largest allegories are epic in scope. A quest forms the narrative thread of both the Greek epic *Odyssey* and the Latin, *Aeneid*, and it is an allegory of the quest for heroic perfection: thus, allegory is aligned with the epic form. Romances, both prose and verse, are inevitably allegorical, although their forms vary in detail with the prevailing cultural ideals of the age. By comparison, the forms of fable and parable are relatively stable—yet even they may play down the moral idea or the mysterious element and emphasize instead the narrative interest, which then results in an elaboration of the form. (Such an elaboration may be seen in a given tale, as told by successive fabulists, such as a fable of the town mouse and the country mouse; with each retelling, the story is absorbed into a new matrix of interpretation.)

Shifts from naïve to sophisticated intent are accompanied by shifts in form. The early authors of fable, following Aesop, wrote in verse; but in the 10th century there appeared collected fables, entitled *Romulus*, written in prose (and books such as this brought down into the medieval and modern era a rich tradition of prose fables). This collection in turn was converted back into elegiac verse. Later masters of fable wrote in verse, but modern favourites—such as Joel Chandler Harris, author of “Uncle Remus” stories, Beatrix Potter, creator of Peter Rabbit, or James Thurber in *Fables for Our Time*—employ their own distinctive prose. Again, while for parables prose narrative may be the norm, they have also been told in verse (as in the emblematic poetry of the 17th-century English Metaphysical poets such as George Herbert, Francis Quarles, and Henry Vaughan).

Loosening the allegorical forms further, some authors have combined prose with verse. Boethius’ *Consolation of Philosophy* (c. AD 524) and Dante’s *The New Life* (c. 1293) interrupt the prose discourse with short poems. Verse and prose then interact to give a new thematic perspective. A related mixing of elements appears in Menippean satire (those writings deriving from the 3rd-century-BC Cynic philosopher Menippus of Gadara), as exemplified in Swift’s *Tale of a Tub*. There a relatively simple allegory of Reformation history (the *Tale* proper) is interrupted by a series of digressions that comment allegorically on the story into which they break.

Even the lyric poem can be adapted to yield allegorical themes and was made to do so, for example, in the visionary and rhapsodic odes written during the high Romantic period after the late 18th century throughout Europe.

The lesson seems to be that every literary genre is adaptable to the allegorical search for multiplicity of meaning.

Diversity of media. In the broadest sense, allegorical procedures are linguistic. Allegory is a manipulation of the language of symbols. Verbally, this mode underwent a major shift in medium along with the shift from oral

to written literature: allegories that had initially been delivered in oral form (Christ’s parables, for example) were written down by scribes and then transcribed by subsequent generations. Much more remarkable transformations, however, take place when the verbal medium is replaced by nonverbal or partially verbal media.

The drama is the chief of such replacements. The enactment of myth in the beginning had close ties with religious ritual, and in the drama of classical Greece both comedy and tragedy, by preserving ritual forms, lean toward allegory. Old Comedy, as represented by the majority of plays by Aristophanes, contains a curious blending of elements—allusions to men of the day, stories suggesting ideas other than the obvious literal sense, religious ceremony, parodies of the graver mysteries, personified abstractions, and stock types of character. Aeschylus’ *Prometheus Bound* uses allegory for tragic ends, while Euripides’ tragedies make a continuous interpretive commentary on the hidden meaning of the basic myths. Allegory is simplified in Roman drama, submitting heroic deeds to the control of the fickle, often malignant goddess Fortuna. Christian symbolism is responsible for the structure of the medieval morality plays, in which human dilemmas are presented through the conflicts of personified abstractions such as the “Virtues” and their “Vice” opponents. The allegory in Renaissance drama is often more atmospheric than structural—though even Shakespeare writes allegorical romances, such as *Cymbeline*, *Pericles*, and *The Winter’s Tale* (and allowed his tragedy of *Coriolanus* to grow out of the “fable of the belly,” which embodies a commonplace of Renaissance political wisdom and is recounted by one of the characters in the play). In 1598 Ben Jonson introduced the comedy of humours, which was dependent on the biological theory that the humours of the body (blood, phlegm, black bile, yellow bile) affect personality; in Jonson’s play *Epicoene, or The Silent Woman* (1609), the character Morose is possessed by the demon of ill humour. Comic allegory of this kind evolved into the Restoration comedy of manners, and through that channel entered modern drama, with Wilde, Shaw, and Pirandello. Ibsen, the master of realistic drama, himself used a free-style allegory in *Peer Gynt*, while the surrealism of modern dramatists such as Ionesco, Genet, and Beckett serves to reinforce the real meaning of their plays.

The degree to which the cinema has been allegorical in its methods has never been surveyed in detail. Any such survey would certainly reveal that a number of basic techniques in film montage builds up multiple layers of meaning. (Animated cartoons, too, continue the tradition of Aesopian fable.)

From time immemorial men have carved religious monuments and have drawn and painted sacred icons. Triumphant arches and chariots have symbolized glory and victory. Religious art makes wide use of allegory, both in its subject matter and in its imagery (such as the cross, the fish, the lamb). Even in poetry there can be an interaction of visual and verbal levels, sometimes achieved by patterning the stanza form. George Herbert’s “Easter Wings,” for instance, has two stanzas set out by the typographer to resemble the shape of a dove’s wings. Such devices belong to the Renaissance tradition of the “emblem,” which combines a motto with a simple symbolic picture (often a woodcut or engraving) and a concise explanation of the picture motto.

While allegory thrives on the visual, it has also been well able to embrace the empty form of pure mathematics. Number symbolism is very old: early Christian systems of cosmology were often based on the number three, referring to the doctrine of the Trinity (and in fact recalling earlier Hebraic and even Hellenic numerology). Musical symbolism has been discovered in the compositions of the 18th-century Baroque composers such as Johann Sebastian Bach. The most evanescent form of allegory, musical imagery and patterns, is also the closest to pure religious vision, since it merges the physical aspects of harmony (based on number) with the sublime and metaphysical effect on its hearers. The final extension of media occurs in the combination of spectacle, drama, dance, and music that is achieved by grand opera, which is at its most

The morality play

Forms that mix prose and poetry

The use of musical symbolism

allegorical in the total artwork of Richard Wagner in the second half of the 19th century. His *Ring* cycle of operas is a complete mythography and allegory, with words and music making two levels of meaning and the whole unified by a type of musical emblem, which Wagner called the leitmotif.

Allegory and cosmology. The allegorical mode has been of major importance in representing the cosmos: the earliest Greek philosophers, for example, speculated on the nature of the universe in allegorical terms; in the Old Testament's oblique interpretation of the universe, too, the world is seen as a symbolic system. The symbolic stories that explain the cosmos are ritualized to ensure that they encode a message. Held together by a system of magical causality, events in allegories are often surrounded by an occult atmosphere of charms, spells, talismans, genies, and magic rites. Science becomes science fiction or a fantastic setting blurs reality so that objects and events become metamorphically unstable. Allegorical fictions are often psychological dramas whose scene is the mind; then their protagonists are personified mental drives. Symbolic climate is most prominent in romance, whose heroic quests project an aura of erotic mysticism, perfect courtesy, and moral fervour that creates a sublime heightening of tone and a picturesque sense of good order.

The cosmic and demonic character of allegorical thinking is most fully reflected in the period of its greatest vogue, the High Middle Ages. During this period poets and priests alike were able to read with increasingly elaborate allegorical technique until their methods perhaps overgrew themselves. A belief had been inherited in the "Great Chain of Being," the Platonic principle of cosmic unity and fullness, according to which the lowest forms of being were linked with the highest in an ascending order. On the basis of this ladderlike conception were built systems of rising transcendence, starting from a material basis and rising to a spiritual pinnacle. The early Church Fathers sometimes used a threefold method of interpreting texts, encompassing literal, moral, and spiritual meanings. This was refined and commonly believed to have achieved its final form in the medieval allegorist's "fourfold theory of interpretation." This method also began every reading with a search for the literal sense of the passage. It moved up to a level of ideal interpretation in general, which was the allegorical level proper. (This was an affirmation that the true Christian believer was right to go beyond literal truth.) Still higher above the literal and the allegorical levels, the reader came to the tropological level, which told him where his moral duty lay. Finally, since Christian thought was apocalyptic and visionary, the fourfold method reached its apogee at the anagogic level, at which the reader was led to meditate on the final cosmic destiny of all Christians and of himself as a Christian hoping for eternal salvation.

While modern scholars have shown that such thinking played its part in the poetry of the Middle Ages and while the Italian poet Dante himself discussed the theological relations between his poems and such a method of exegesis, the main arena for the extreme elaboration of this allegory was in the discussion and the teaching of sacred Scriptures. As such, the fourfold method is of highest import, and it should be observed that it did not need to be applied in a rigid four-stage way. It could be reduced, and commonly was reduced, to a two-stage method of interpretation. Then the reader sought simply a literal and a spiritual meaning. But it could also be expanded. The passion for numerology, combined with the inner drive of allegory toward infinite extension, led to a proliferation of levels. If four levels were good, then five or eight or nine might be better.

HISTORICAL DEVELOPMENT IN WESTERN CULTURE

Fable. The origins of fable are lost in the mists of time. Fables appear independently in ancient Indian and Mediterranean cultures. The Western tradition begins effectively with Aesop (6th century BC), of whom little or nothing is known for certain; but before him the Greek poet Hesiod (8th century BC) recounts the fable of the hawk and the nightingale, while fragments of similar tales

survive in Archilochus, the 7th-century-BC warrior-poet. Within 100 years of the first Aesopian inventions the name of Aesop was firmly identified with the genre, as if he, not a collective folk, were its originator. Like the Greek philosopher Socrates, Aesop was reputed to have been ugly but wise. Legend connected him with the island of Samos; the historian Herodotus believed him to have been a slave.

Modern editions list approximately 200 "Aesop" fables, but there is no way of knowing who invented which tales or what their original occasions might have been. Aesop had already receded into legend when Demetrius of Phaleron, a rhetorician, compiled an edition of Aesop's fables in the 4th century BC. The poetic resources of the form developed slowly. A versified Latin collection made by Phaedrus, a freed slave in the house of the Roman emperor Augustus, included fables invented by the poet, along with the traditional favourites, which he retold with many elaborations and considerable grace. (Phaedrus may also have been the first to write topically allusive fables, satirizing Roman politics.) A similar extension of range marks the work of the Hellenized Roman Babrius, writing in the 2nd century AD. Among the classical authors who advanced upon Aesopian formulas may be named the Roman poet Horace, the Greek biographer Plutarch, and the great satirist Lucian of Samosata.

Beast epic. In the Middle Ages, along with every other type of allegory, fable flourished. Toward the end of the 12th century, Marie de France made a collection of over 100 tales, mingling beast fables with stories of Greek and Roman worthies. In another compilation, Christine de Pisan's Othéa manuscript illuminations provide keys to the interpretation of the stories and support the appended moral tag line. Expanded, the form of the fable could grow into what is called the beast epic, a lengthy, episodic animal story, replete with hero, villain, victim, and endless epic endeavour. (One motive for thus enlarging upon fable was the desire to parody epic grandeur: the beast epic mocks its own genre.) Most famous of these works is a 12th-century collection of related satirical tales called *Roman de Renart*, whose hero is a fox, symbolizing cunning man. The *Roman* includes the story of the fox and Chantecler (Chanticleer), a cock, a tale soon afterward told in German, Dutch, and English versions (in *The Canterbury Tales*, Geoffrey Chaucer took it as the basis for his "Nun's Priest's Tale"). Soon the *Roman* had achieved universal favour throughout Europe. The Renaissance poet Edmund Spenser also made use of this kind of material; in his "Mother Hubbard's Tale," published in 1591, a fox and an ape go off to visit the court, only to discover that life is no better there than in the provinces. More sage and serious, John Dryden's poem of *The Hind and Panther* (1687) revived the beast epic as a framework for theological debate. Bernard de Mandeville's *Fable of the Bees* (first published 1705 as *The Grumbling Hive, or Naves Turn'd Honest*) illustrated the rapacious nature of humans in society through the age-old metaphor of the kingdom of the bees. In modern times, children's literature has made use of animal fable but often trivialized it. But the form has been taken seriously, as, for example, by the political satirist George Orwell, who, in his novel *Animal Farm* (1945), used it to attack Stalinist Communism.

Influence of Jean de La Fontaine. The fable has normally been of limited length, however, and the form reached its zenith in 17th-century France, at the court of Louis XIV, especially in the work of Jean de La Fontaine. He published his *Fables* in two segments: the first, his initial volume of 1668, and the second, an accretion of "Books" of fables appearing over the next 25 years. The 1668 *Fables* follow the Aesopian pattern, but the later ones branch out to satirize the court, the bureaucrats attending it, the church, the rising bourgeoisie—indeed the whole human scene. La Fontaine's great theme was the folly of human vanity. He was a skeptic, not unkind but full of the sense of human frailty and ambition. His satiric themes permitted him an enlargement of poetic diction; he could be eloquent in mocking eloquence or in contrast use a severely simple style. (His range of tone and style has been admirably reflected in a version of his works made by a

The "fourfold theory of interpretation"

Tales of the cunning fox

20th-century American poet, Marianne Moore.) La Fontaine's example gave new impetus to the genre throughout Europe, and during the Romantic period a vogue for Aesopian fable spread to Russia, where its great practitioner was Ivan Andreyevich Krylov. The 19th century saw the rise of literature written specifically for children, in whom fable found a new audience. Among the most celebrated authors who wrote for them are Lewis Carroll, Charles Kingsley, Rudyard Kipling, Kenneth Grahame, Hilaire Belloc, and Beatrix Potter. There is no clear division between such authors and the "adult" fabulist, such as Hans Christian Andersen, Lewis Carroll, Oscar Wilde, Saint-Exupéry, or J.R.R. Tolkien. In the 20th century there are the outstanding *Fables for Our Time*, written by James Thurber and apparently directed toward an adult audience (although a sardonic parent might well read the *Fables* to his children).

Parable. In the West, the conventions of parable were largely established by the teachings of Christ. The New Testament records a sufficient number of his parables, with their occasions, to show that to some extent his disciples were chosen as his initiates and followers because they "had ears to hear" the true meaning of his parables. (It has already been noted that the parable can be fully understood only by an elite, made up of those who can decipher its inner core of truth.) Despite a bias toward simplicity and away from rhetorical elaboration, the parable loses little in the way of allegorical richness: the speaker can exploit an enigmatic brevity that is akin to the style of presenting a complex riddle. Parable is thus an immensely useful preaching device; while theologians in the period of the early Christian Church were developing glosses on Christ's enigmatic stories, preachers were inventing their own to drive home straightforward lessons in good Christian conduct. For centuries, therefore, the model of parable that had been laid down by Christ flourished on Sundays in churches all over the Western world. Pious tales were collected in handbooks: the *Gesta Romanorum*, the *Alphabet of Thales*, the *Book of the Knight of La Tour Landry*, and many more. Infinitely varied in subject matter, these exemplary tales used a plain but lively style, presenting stories of magicians, necromancers, prophets, chivalrous knights and ladies, great emperors—a combination bound to appeal to congregations, if not to theologians. An important offshoot of the parable and exemplary tale was the saint's life. Here, too, massive compilations were possible; the most celebrated was *The Golden Legend* of the 15th century, which included approximately 200 stories of saintly virtue and martyrdom.

20th-century parables

As long as preaching remained a major religious activity, the tradition of parable preserved its strong didactic strain. Its more paradoxical aspect gained renewed lustre in theological and literary spheres when the 19th-century Danish philosopher Soren Kierkegaard began to use parables in his treatises on Christian faith and action. In *Fear and Trembling* he retold the story of Abraham and Isaac; in *Repetition* he treated episodes in his own life in the manner of parable. Such usage led to strange new literary forms of discourse, and his writing influenced, among others, the Austro-Czech novelist Franz Kafka and the French "absurdist" philosopher, novelist, and playwright Albert Camus. Kafka's parables, full of doubt and anxiety, mediate on the infinite chasm between man and God and on the intermediate role played by the law. His vision, powerfully expressed in parables of novel length (*The Castle*, *The Trial*, *Amerika*), is one of the most enigmatic in modern literature.

Allegory. The early history of Western allegory is intricate and encompasses an interplay between the two prevailing world views—the Hellenic and the Hebraic-Christian—as theologians and philosophers attempted to extract a higher meaning from these two bodies of traditional myth.

In terms of allegory, the Greco-Roman and Hebraic-Christian cultures both have a common starting point: a creation myth. The Old Testament's book of Genesis roughly parallels the story of the creation as told by the Greek poet Hesiod in his *Theogony* (and the later Roman version of the same event given in Ovid's *Metamorpho-*

ses). The two traditions thus start with an adequate source of cosmic imagery, and both envisage a universe full of mysterious signs and symbolic strata. But thereafter the two cultures diverge. This is most apparent in the way that the style of the body of poetry attributed to Homer—the ancient Greek "Bible"—differs from the Old Testament narrative. The Greek poet presented his heroes against an articulated narrative scene, a context full enough for the listener (and, later, the reader) to ignore secondary levels of significance. By contrast, the Jewish authors of the Old Testament generally emptied the narrative foreground, leaving the reader to fill the scenic vacuum with a deepening, thickening allegorical interpretation.

Old Testament. The Old Testament, including its prophetic books, has a core of historical record focussing on the trials of the tribes of Israel. In their own view an elect nation, the Israelites believe their history spells out a providential design. The prophets understand the earliest texts, Genesis and Exodus, in terms of this providential scheme. Hebraic texts are interpreted as typological: that is, they view serious myth as a theoretical history in which all events are types—portents, foreshadowing the destiny of the chosen people. Christian exegesis (the critical interpretation of Scripture) inherits the same approach.

Typological allegory looks for hidden meaning in the lives of actual men who, as types or figures of later historical persons, serve a prophetic function by prefiguring those later persons. Adam, for example (regarded as a historical person), is thought to prefigure Christ in his human aspect. Joshua to prefigure the victorious militant Christ. This critical approach to Scripture is helped by the fact of monotheism, which makes it easier to detect the workings of a divine plan. The splendours of nature hymned in the Psalms provide a gloss upon the "glory of God." The Law (the Torah) structures the social aspect of sacred history and, as reformulated by Christ, provides the chief link between Old and New Testaments. Christ appeals to the authority of "the Law and the Prophets" but assumes the ultimate prophetic role himself, creating the New Law and the New Covenant—or Testament—with the same one God of old.

Typological allegory

The Greeks. Hellenic tradition after Homer stands in sharp contrast to this concentration on the fulfilling of a divine plan. The analytic, essentially scientific histories of Herodotus and Thucydides precluded much confident belief in visionary providence. The Greeks rather believed history to be structured in cycles, as distinct from the more purposive linearity of Hebraic historicism.

Nevertheless, allegory did find a place in the Hellenic world. Its main arena was in philosophic speculation, centring on the interpretation of Homer. Some philosophers attacked and others defended the Homeric mythology. A pious defense argued that the stories—about the monstrous love affairs of the supreme god Zeus, quarrels of the other Olympian gods, scurrility of the heroes, and the like—implied something beyond their literal sense. The defense sometimes took a scientific, physical form: in this case, Homeric turmoil was seen as reflecting the conflict between the elements. Or Homer was moralized: the goddess Pallas Athene, for example, who in physical allegory stood for the ether, in moral allegory was taken to represent reflective wisdom because she was born out of the forehead of her father, Zeus. Moral and physical interpretation is often intermingled.

Plato, the Idealist philosopher, occupies a central position with regard to Greek allegory. His own myths imply that our world is a mere shadow of the ideal and eternal world of forms (the Platonic ideas), which has real, independent existence, and that the true philosopher must therefore be an allegorist in reverse. He must regard phenomena—things and events—as a text to be interpreted upward, giving them final value only insofar as they reveal, however obscurely, their ideal reality in the world of forms. Using this inverted allegorical mode, Plato attacked Homeric narrative, whose beauty beguiles men into looking away from the truly philosophic life. Plato went further. He attacked other fashionable philosophic allegorists because they did not lead up to the reality but limited speculation to the sphere of moral and physical necessity. Platonic

allegory envisaged the system of the universe as an ascending ladder of forms, a "Great Chain of Being," and was summed up in terms of myth in his *Timaeus*. Plato and Platonic thought became, through the influence of this and other texts on Plotinus (died 269/270) and through him on Porphyry (died c. 304), a pagan mainstay of later Christian allegory. Medieval translations of Dionysius the Areopagite (before 6th century AD) were equally influential descendants of Platonic vision.

A second and equally influential Hellenic tradition of allegory was created by the Stoic philosophers, who held that the local gods of the Mediterranean peoples were signs of a divinely ordered natural destiny. Stoic allegory thus emphasized the role of fate, which, because all men were subject to it, could become a common bond between peoples of different nations. A later aspect of moral exegesis in the Stoic manner was the notion that myths of the gods really represent, in elevated form, the actions of great men. In the 2nd century BC, under Stoic influence, the Sicilian writer Euhemerus argued that theology had an earthly source. His allegory of history was the converse of Hebraic typology—which found the origin of the divine in the omnipotence of the One God—for Euhemerus found the origin of mythological gods in human kings and heroes, divinized by their peoples. His theories enjoyed at least an aesthetic revival during the Renaissance.

Blending of rival systems: the Middle Ages. At the time of the birth of Christ, ideological conditions within the Mediterranean world accelerated the mingling of Hellenic and Hebraic traditions. Philo Judaeus laid the groundwork; Clement of Alexandria and Origen followed him. The craft of allegorical syncretism—that is, making rival systems accommodate one another through the transformation of their disparate elements—was already a developed art by the time St. Paul and the author of the Gospel According to St. John wove the complex strands of the Hebraic-Christian synthesis. Over centuries of quarrelling, the timeless philosophy of the Greek allegorists was accommodated to the time-laden typology of the Hebrew prophets and their Christian successors and at length achieved a hybrid unity that permitted great allegories of Western Christendom to be written.

As a hybrid method, allegory could draw on two archetypal story lines: the war and the quest of Homer's *Iliad* and *Odyssey*, which was paralleled by the struggles and wanderings of the children of Israel. Throughout the Middle Ages the figure of the wandering Aeneas (who, in the second half of Virgil's Latin epic, *Aeneid*, fought bloody battles) was seen as a type in a system of hidden Christianity. Virgil's fourth *Eclogue*, a prophetic vision of the birth of a child who would usher in the "golden age," was read as a prophecy of the birth of Christ. Seen by many Christian commentators as the ideal allegorist, Virgil himself was hailed as a proto-Christian prophet. The blending of rival systems of allegory from widely assorted cultures became the rule for later allegory. Adapting the Latin writer Apuleius' fable of Cupid and Psyche, Edmund Spenser combined its elements with ancient Middle Eastern lore, Egyptian wisdom, and dashes of Old Testament critical interpretation to convert the enclosed garden of the biblical Song of Solomon into the gardens of Adonis in *The Faerie Queene*, Book III. The pagan gods survived unharmed throughout the Middle Ages if wearing Christian costumes, because Christians were taught that pagan worthies could be read as figures of Christian rulers. The labours of Hercules, for instance, stood for the wanderings and trials of all Christian men; the Hellenic theme of heroic warfare took a Christianized form, available to allegory, when in the 4th century the poet and hymn writer Prudentius internalized war as the inner struggle of Christian man, suspended between virtue and vice. For complete triumph in explaining the significance of the world, Christianity needed one further element: a world-historical theory large enough to contain all other theories of meaning. This it found in the belief that God was the author of the world. His creation wrote the world. The world, read as a text, provided a platform for transforming the piecemeal, postclassical syncretism into some semblance of order. Firmly established in the West,

Christianity, for all its strains of discord, slowly achieved a measure of coherence. St. Thomas Aquinas could write its *Summa*. Theocentric, authoritarian, spiritualist, and word oriented, the medieval model of allegory lent itself to the creation of the most wonderful of all allegorical poems. Dante's *Divine Comedy*, completed shortly before his death in 1321.

Before this could happen, however, the Christian world view had been subjected to an important pressure during the 12th century. It may be called the pressure to externalize. Alanus de Insulis (Alain de Lille), Bernard of Sylvestris, John of Salisbury, and other forerunners of the movement known as European Humanism "discovered" nature. Delighting in the wonders of God's cosmic text, they brought theological speculation down to earth. Romances of love and chivalry placed heroes and heroines against the freshness of spring. Everywhere nature shone, sparkling with the beauty of earthly life. The externalization and naturalizing of Christian belief flowers most obviously in *The Romance of the Rose*, begun in the 13th century by Guillaume de Lorris and completed, in vastly complicated form, by Jean de Meung. The *Romance* personifies the experiences of courtly love, recounting the pursuit of an ideal lady by an ideal knight, set in an enclosed garden and castle, which permits Guillaume to dwell on the beauty of nature. With Jean de Meung the interest in nature is made explicit, and the poem ends in a series of lengthy digressive discourses, several of them spoken by Dame Nature herself. In medieval English poetry this same love of spring and seasonal pleasures is apparent everywhere—certainly in the poems of Geoffrey Chaucer, who, besides creating several allegories of his own, translated *The Romance of the Rose* into English.

Dante's *Divine Comedy* has physical immediacy and contains an immense amount of historical detail. He anchors his poem in a real world, accepting Christian typology as historical fact and adopting an ordered system of cosmology (based on the number three, proceeding from the Trinity). Dante's passion for numerology does not, however, block a closeness to nature that had perhaps not been equalled in poetry since Homer. He enfold classical thought into his epic by making Virgil one of its main protagonists—again to prefigure Christian heroism. Perhaps only William Langland, the author of *The Vision of Piers Plowman*, could be said to rival Dante's cosmic range. *Piers Plowman* is a simpler apocalyptic vision than the *Comedy*, but it has an existential immediacy, arising from its concern for the poor, which gives it great natural power.

Renaissance. Romance and romantic forms provide the main vehicle for the entrance of allegory into the literature of the Renaissance period. The old Arthurian legends carry a new sophistication and polish in the epics of the Italians Boiardo, Ariosto, and Tasso and in the work of Edmund Spenser. By interlacing several simultaneous stories in one larger narrative, the literary technique known as *entrelacement* allowed digression—yet kept an ebbing, flowing kind of unity—while presenting opportunities for moral and ironic commentary. But although the forms and themes of romance were medieval in origin, the new age was forced to accommodate altered values. The Middle Ages had externalized the Christian model; the Renaissance now internalized it, largely by emphasizing the centrality of human understanding. This process of internalization had begun slowly. In rough outline it can be discerned in the belief that biological humours affected personality, in the adaptations of Platonic idealism from which arose a new emphasis on imagination, in the rise of an introspective, soliloquizing drama in England. It can further be discerned in the gradual adoption of more self-conscious theories of being: Shakespeare's Hamlet, finding himself by thinking out his situation, prefigures the first modern philosopher, René Descartes, whose starting point for argument was "I think, therefore I am." Christopher Marlowe's characterization of Dr. Faustus epitomizes the new age. Pursuing power in the form of knowledge, he is led to discover the demons of allegory within himself. He is an essential figure for later European literature, archetypal in Germany for both Johann Wolfgang von Goethe and Thomas Mann and influential everywhere.

Stoic
allegory

The
importance
of the
Christian
world view
in allegory

Milton's
allegory:
*Paradise
Lost*

Modern period. With the Baroque and Neoclassical periods, allegory began to turn away from cosmology and toward rhetorical ambiguity. John Milton allegorized sin and death in his epic poem *Paradise Lost*, but allegory for him seems chiefly to lie in the ambiguous diction and syntax employed in the poem. Instead of flashing allegorical emblems before the reader, Milton generates a questioning attitude that searches out allegory more as a mysterious form than a visible content. His central allegorical theme is perhaps the analogy he draws between poetry, music, and ideas of cosmic order. This theme, which generates allegory at once, recurs in later English poetry right up to modern times with T.S. Eliot's *Four Quartets*.

The social and religious attitudes of the 18th-century age of Enlightenment could be expressed coolly and without ambiguity—and thus there was little need for spiritual allegory in the period's literature. Oblique symbolism was used mainly for satirical purposes. John Dryden and Alexander Pope were masters of verse satire, Jonathan Swift of prose satire. Voltaire and the French writers of the Enlightenment similarly employed a wit whose aim was to cast doubt on inherited pieties and attitudes. A new vogue for the encyclopaedia allowed a close, critical commentary on the ancient myths, but the criticism was rationalist and opposed to demonology. Under such conditions the allegorical mode might have dried up entirely. Yet the new Romantic age of the late 18th and early 19th centuries revived the old cosmologies once more, and poetic forms quickly reflected the change, with the Romantic poets and their precursors (Blake, William Collins, Edward Young, Thomas Gray, and others) managing to reinstate the high destiny of the allegorical imagination. The Romantics went back to nature. Poets took note of exactly what they saw when they went out walking, and their awareness of nature and its manifestations found its way into their poetry. Appropriate poetic forms for expressing this sensibility tended to be open, rhapsodic, and autobiographical—qualities notably present in William Wordsworth and in Samuel Taylor Coleridge, for example. Percy Bysshe Shelley is the most strikingly allegorical of English Romantics; he not only followed the Platonic tradition of Spenser and the Renaissance—with ode, elegy, and brief romance—but he also invented forms of his own, such as *Epipsychidion*, a rhapsodic meditation, and he was working on a great Dantesque vision, *The Triumph of Life*, when he died. Visionary masterpieces came from Germany, where Novalis and Friedrich Hölderlin hymned the powers of nature in odes of mythic overtone and resonance. French Romanticism, merging gradually with the theory and practice of the Symbolist movement (dealing in impressions and intuitions rather than in descriptions), in turn followed the same path. The pantheist cosmology of Victor Hugo, the central writer of the somewhat delayed French Romantic movement, created an allegory of occult forces and demonic hero worship. It is fair to say that, in its most flexible and visionary forms, allegory flourished throughout the Romantic period.

There also developed a novelistic mode of allegory by which prose authors brought fate, necessity, the demonic, and the cosmological into their narratives. Émile Zola used a theory of genetics, Charles Dickens the idea of ecological doom, Leo Tolstoy the belief in historical destiny, and Fyodor Dostoyevsky the fatalism of madness and neurosis. Nikolay Gogol revived the art of the grotesque, picturing absurdities in the scene of tsarist Russia. Even the arch-naturalist playwright, Anton Chekhov, made an emblem of the cherry orchard and the sea gull in his plays of those titles. However its dates are established, the modern period is exceedingly complex in its mythmaking. Psychoanalytic theory has been both a critical and a creative resource; modern allegory has remained internalized in the Renaissance tradition. But Marxist social realism has kept to the externals of dialectical materialism, though without notable aesthetic success. In the free play of American letters, where Nathaniel Hawthorne, Herman Melville, Edgar Allan Poe, and Henry James (particularly in his later novels) had essayed an allegorical mode, the future of its use is uncertain. T.S. Eliot's enigmatic style in a long poem, "Ash Wednesday," may be related to his

Complexity of
myth-
making in
the 20th
century

search for a Dantesque dramatic style, for which he also tried in plays, most obviously *Murder in the Cathedral* (a morality) and *The Cocktail Party* (a philosophic farce). More clearly popular authors such as George Orwell and William Golding have used the most familiar allegorical conventions. D.H. Lawrence shaped novels such as *The Plumed Serpent* to project a thematic, cultural polemic. W.H. Auden's operatic librettos reflect once more the allegorical potential of this mixture of media.

Modern allegory has in fact no set pattern, or model, although Surrealism has provided a dominant style of discontinuous fragmentary expression. The only rule seems to be that there is no rule. Science fiction, an ancient field dating back at least to the earliest philosophers of Greece, has set no limits on the speculations it will entertain. The allegorical author now even questions the allegorical process itself, criticizing the very notions of cosmos, demon, and magic. It may be that modern allegory has completed a vast circle begun by the first conflict between ways of interpreting myth, as revealed in Homer and the Hebraic prophets.

(A.S.F.)

ALLEGORICAL LITERATURE IN THE EAST

India. Fables appeared early in India, but it is impossible to determine whether they are older or later than the Greek. Undoubtedly there was mutual influence from very early times, for indirect contacts between Greece and India (by trade routes) had existed long before the time of Alexander the Great. In the form in which they are now known the Greek fables are the older, but this may be an accident of transmission.

The fable was apparently first used in India as a vehicle of Buddhist instruction. Some of the *Jātakas*, birth stories of the Buddha, which relate some of his experiences in previous animal incarnations, resemble Greek fables and are used to point a moral. They may date from as far back as the 5th century BC, though the written records are much later. The most important compilation is the *Pañcatantra*, a Sanskrit collection of beast fables. The original has not survived, but it has been transmitted (via a lost Pahlavi version) as the mid-8th-century Arabic *Kalīlah wa Dimnah*. Kalilah and Dimnah are two jackals, counselors to the lion king, and the work is a frame story containing numerous fables designed to teach political wisdom or cunning. From the Arabic this was translated into many languages, including Hebrew, which version John of Capua used to make a Latin version in the 13th century. This, the *Directorium humanae vitae* ("Guide for Human Life"), was the chief means by which oriental fables became current in Europe. In the fables of Bidpai, animals act as men in animal form, and little attention is paid to their supposed animal characteristics. It is in this respect that they differ most from the fables of Aesop, in which animals behave as animals.

(Ed.)

China. Chinese philosophers from the Ch'in dynasty (221–206 BC) onward often used extended metaphors (from which fable is the logical development) to make their points. This is believed to reflect the fact that, as "realistic" thinkers, the Chinese generally did not favour more abstract argument. Thus simple allegory helped to stimulate audience interest and to increase the force of an argument. A century earlier, Mencius, a Confucian philosopher, had used the following little allegory in illustrating his theory that an effort has to be made if man's natural goodness is to be recovered: "A man will begin searching when his dog or chicken is missing; but he does not go searching for the good character he was born with after it is lost. Is this not regrettable?" The same writer also used a parable to bring home his point that mental training could not be hurried, but was a gradual process: "A man in Sung sowed seeds in a field. The seedlings grew so slowly, however, that one day he took a walk through the field pulling at each one of the seedlings. On returning home he announced that he was exhausted, but that he had helped the seedlings' growth. His son, hurrying to the field, found the seedlings dead."

Tales such as this were often borrowed from folklore,

but others were probably original creations, including a striking story that opens the *Chuang-tzu*, a summa of Taoist thought. It makes the point that ordinary people frequently deplore the actions of a man of genius because they are unable to understand his vision, which is not answerable to the laws of "common sense": "A giant fish, living at the northern end of the world, transformed itself into a bird so that it could make the arduous flight to the southernmost sea. Smaller birds, measuring his ambition against their own capabilities, laughed at the impossibility of it."

But the full development of fable, as it is understood in the West, was hindered by the fact that Chinese ways of thinking prohibited them from accepting the notion of animals that thought and behaved as humans. Actual events from the past were thought to be more instructive than fictitious stories, and this led to the development of a large body of legendary tales and supernatural stories. Between the 4th and 6th centuries, however, Chinese Buddhists adapted fables from Buddhist India in a work known as *Po-Yü ching*, and they also began to make use of traditional Chinese stories that could further understanding of Buddhist doctrines. (Na.Mo.)

Japan. In Japan, the *Koji-ki* (712; "Records of Ancient Matters") and the *Nihon-shoki* (8th century, "Chronicles of Japan"), both of them official histories of Japan, were studded with fables, many of them on the theme of a small intelligent animal getting the better of a large stupid one. The same is true of the *fudoki* (local gazetteers dating from 713 and later). The form reached its height in the Kamakura period (1192–1333). Toward the end of the Muromachi period (1338–1573) Jesuit missionaries introduced the fables of Aesop to Japan, and the influence of these can be traced in stories written between then and the 19th century. (T.Io.)

Ballad

The ballad is a short narrative folk song whose distinctive style crystallized in Europe in the late Middle Ages and persists to the present day in communities where literacy, urban contacts, and mass media have not yet affected the habit of folk singing. France, Denmark, Germany, Russia, Greece, and Spain, as well as England and Scotland, possess impressive ballad collections. At least one-third of the 300 extant English and Scottish ballads have counterparts in one or several of these continental balladries, particularly those of Scandinavia. In no two language areas, however, are the formal characteristics of the ballad identical. For example, British and American ballads are invariably rhymed and strophic (*i.e.*, divided into stanzas); the Russian ballads known as *byliny* and almost all Balkan ballads are unrhymed and unstrophic; and, though the *romances* of Spain, as their ballads are called, and the Danish *viser* are alike in using assonance instead of rhyme, the Spanish ballads are generally unstrophic while the Danish are strophic, parcelled into either quatrains or couplets.

ELEMENTS

Narrative basis. Typically, the folk ballad tells a compact little story that begins abruptly at the moment when the narrative has turned decisively toward its catastrophe or resolution. Focussing on a single, climactic situation, the ballad leaves the inception of the conflict and the setting to be inferred or sketches them in hurriedly. Characterization is minimal, the characters revealing themselves in their actions or speeches; overt moral comment on the characters' behaviour is suppressed and their motivation seldom explicitly detailed. Whatever description occurs in ballads is brief and conventional: transitions between scenes are abrupt and time shifts are only vaguely indicated; crucial events and emotions are conveyed in crisp, poignant dialogue. In short, the ballad method of narration is directed toward achieving a bold, sensational, dramatic effect with purposeful starkness and abruptness. But despite the rigid economy of ballad narratives, a repertory of rhetorical devices is employed for prolonging highly charged moments in the story and thus thickening

the emotional atmosphere. In the most famous of such devices, incremental repetition, a phrase or stanza is repeated several times with a slight but significant substitution at the same critical point. Suspense accumulates with each substitution, until at last the final and revelatory substitution bursts the pattern, achieving a climax and with it a release of powerful tensions. The following stanza is a typical example:

Then out and came the thick, thick, blood,
Then out and came the thin,
Then out and came the bonny heart's blood,
Where all the life lay in.

Oral transmission. Since ballads thrive among unlettered people and are freshly created from memory at each separate performance, they are subject to constant variation in both text and tune. Where tradition is healthy and not highly influenced by literary or other outside cultural influences, these variations keep the ballad alive by gradually bringing it into line with the style of life, beliefs, and emotional needs of the immediate folk audience. Ballad tradition, however, like all folk arts, is basically conservative, a trait that explains the references in several ballads to obsolete implements and customs, as well as the appearance of words and phrases that are so badly garbled as to indicate that the singer does not understand their meaning though he takes pleasure in their sound and respects their traditional right to a place in his version of the song. The new versions of ballads that arise as the result of cumulative variations are no less authentic than their antecedents. A poem is fixed in its final form when published, but the printed or taped record of a ballad is representative only of its appearance in one place, in one line of tradition, and at one moment in its protean history. The first record of a ballad is not its original form but merely its earliest recorded form, and the recording of a ballad does not inhibit tradition from varying it subsequently into other shapes, because tradition preserves by re-creating rather than by exact reproduction.

COMPOSITION

Theories. How ballads are composed and set afloat in tradition has been the subject of bitter quarrels among scholars. The so-called communal school, which was led by two American scholars F.B. Gummere (1855–1919) and G.L. Kittredge (1860–1941), argued at first that ballads were composed collectively during the excitement of dance and song festivals. Under attack the communalists retreated to the position that although none of the extant ballads had been communally composed, the prototypical ballads that determined the style of the ballads had originated in this communal fashion. Their opponents were the individualists, who included the British men of letters W.J. Courthope (1842–1917) and Andrew Lang (1844–1912) and the American linguist Louise Pound (1872–1958). They held that each ballad was the work of an individual composer, who was not necessarily a folk singer, tradition serving simply as the vehicle for the oral perpetuation of the creation. According to the widely accepted communal re-creation theory, put forward by the American collector Phillips Barry (1880–1937) and the scholar G.H. Gerould (1877–1953), the ballad is conceded to be an individual composition originally. This fact is considered of little importance because the singer is not expressing himself individually, but serving as the deputy of the public voice, and because a ballad does not become a ballad until it has been accepted by the folk community and been remolded by the inevitable variations of tradition into a communal product. Ballads have also been thought to derive from art songs, intended for sophisticated audiences, which happened to filter down to a folk level and become folk song. This view, though plausible in the case of certain folk lyrics, is inapplicable to the ballads, for if the ballads were simply miscellaneous castoffs, it would not be possible to discern so clearly in them a style that is unlike anything in sophisticated verse.

Technique and form. Ballads are normally composed in two kinds of stanzas; the first consists of a couplet of lines each with four stressed syllables, and with an interwoven refrain:

Rhetorical
devices

But it would have made your heart right sair.

With a hey ho and a lillie gay
To see the bridegroom rive his haire.
As the primrose spreads so sweetly

the second a stanza of alternating lines of four stresses and three stresses, the second and fourth lines rhyming:

There lived a wife at Usher's Well,
And a wealthy wife was she;
She had three stout and stalwart sons,
And sent them o'er the sea.

Reference to the tunes show that the three-stress lines actually end in an implied fourth stress to match the pause in the musical phrase at these points. The interwoven refrain is a concession to the musical dimension of the ballad; it may be a set of nonsense syllables (Dillum down dillum, Fa la la la) or irrelevant rigmaroles of flowers or herbs. A few ballads have stanza-length burdens interspersed between the narrative stanzas, a technique borrowed from the medieval carols. The lyrical and incantatory effect of refrains during the ballad performance is very appealing, but in cold print they often look ridiculous, which is perhaps why early collectors failed to note them. In the first example above, it will be noted that the gaiety of the refrain is at odds with the mood of the meaningful lines. Not infrequently the ballad stanza satisfies the music's insistence on lyrical flourishes by repeating textual phrases and lines:

So he ordered the grave to be opened wide,
And the shroud to be turned down;
And there he kissed her clay cold lips
Till the tears came trickling down, down, down,
Till the tears came trickling down

Kinds of
repetition
used in
ballads

The refrain is just one of the many kinds of repetition employed in ballads. Incremental repetition, already discussed, is the structural principle on which whole ballads ("The Maid Freed from the Gallows," "Lord Randal") are organized, and many other ballads contain long exchanges of similarly patterned phrases building cumulatively toward the denouement:

"Oh what will you leave to your father dear?"
"The silver-shod steed that brought me here."
"What will you leave to your mother dear?"
"My velvet pall and my silken gear."
"What will you leave to your brother John?"
"The gallows-tree to hang him on."

Any compressed narrative of sensational happenings told at a high pitch of feeling is bound to repeat words and phrases in order to accommodate the emotion that cannot be exhausted in one saying, a tendency that accounts for such stanzas as:

Then He says to His mother, "Oh: the withy [willow], oh:
the withy,
The bitter withy that causes me to smart, to
smart,
Oh: the withy, it shall be the very first tree
That perishes at the heart."

Much repetition in ballads is mnemonic as well as dramatic. Since ballads are performed orally, the hearer cannot turn back a page to recover a vital detail that slipped by in a moment of inattention. Crucial facts in narrative, therefore, are incised in the memory by skillful repetition; instructions given in a speech are exactly repeated when the singer reports the complying action; answers follow the form of the questions that elicited them.

Conven-
tional
imagery

The exigencies of oral performance also account for the conventional stereotyped imagery of the ballads. For unlike the poet, who reaches for the individualistic, arresting figure of speech, the ballad singer seldom ventures beyond a limited stock of images and descriptive adjectives. Knights are always gallant, swords royal, water wan, and ladies gay. Whatever is red is as red as blood, roses, coral, rubies, or cherries; white is stereotyped as snow white, lily white, or milk white. Such conventions fall into place almost by reflex action, easing the strain on the singer's memory and allowing him to give his full attention to the manipulation of the story. The resulting bareness of verbal texture, however, is more than compensated for by the dramatic rhetoric through which the narrative is projected. In any case, complex syntax and richness of

language are forbidden to texts meant to be sung, for music engages too much of the hearer's attention for him to untangle an ambitious construction or relish an original image. Originality indeed, like anything else that exalts the singer, violates ballad decorum, which insists that the singer remain impersonal.

Music. A ballad is not technically a ballad unless it is sung; but though tunes and texts are dynamically interdependent, it is not unusual to find the same version of a ballad being sung to a variety of tunes of suitable rhythm and metre or to find the same tune being used for several different ballads. And just as there are clusters of versions for most ballads, so a given ballad may have associated with it a family of tunes whose members appear to be versions of a single prototypical form.

Ballad tunes are based on the modes rather than on the diatonic and chromatic scales that are used in modern music. Where chromaticism is detected in American folk music, the inflected tones are derived from black folk practice or from learned music. Of the six modes, the preponderance of folk tunes are Ionian, Dorian, or Mixolydian; Lydian and Phrygian tunes are rare. The folk music least affected by sophisticated conditioning does not avail itself of the full seven tones that compose each of the modal scales. Instead, it exhibits gapped scales, omitting either one of the tones (hexatonic) or two of them (pentatonic). Modulation sometimes occurs in a ballad from one mode to an adjacent mode.

Most tunes consist of 16 bars with duple rhythm, or two beats per measure, prevailing slightly over triple rhythm. The tune, commensurate with the ballad stanza, is repeated as many times as there are stanzas. Unlike the "through-composed" art song, where the music is given nuances to correspond to the varying emotional colour of the content, the folk song affords little opportunity to inflect the contours of the melody. This limitation partly explains the impassive style of folk singing. Musical variation, however, is hardly less frequent than textual variation; indeed, it is almost impossible for a singer to perform a ballad exactly the same way twice. The stablest part of the tune occurs at the mid-cadence (the end of the second text line) and the final cadence (the end of the fourth line). The third phrase of the tune, corresponding to the third line of the stanza, proves statistically the most variable. Significantly, these notes happen to coincide with the rhyming words. The last note of the tune, the point of resolution and final repose, usually falls on the fundamental tone (*i.e.*, keynote) of the scale; the mid-cadence falling normally a perfect fifth above the tonic or a perfect fourth below it. To make for singability, the intervals in the melodic progression seldom involve more than three degrees. And since the singer performs solo or plays the accompanying instrument himself, he need not keep rigidly to set duration or stress but may introduce grace notes to accommodate hypermetric syllables and lengthen notes for emphasis.

Impassive
style of
folk
singing

TYPES OF BALLADRY

The traditional folk ballad, sometimes called the Child ballad in deference to Francis Child, the scholar who compiled the definitive English collection, is the standard kind of folk ballad in English and is the type of balladry that this section is mainly concerned with. But there are peripheral kinds of ballads that must also be noticed in order to give a survey of balladry.

Minstrel ballad. Minstrels, the professional entertainers of nobles, squires, rich burghers, and clerics until the 17th century, should properly have had nothing to do with folk ballads, the self-created entertainment of the peasantry. Minstrels sometimes, however, affected the manner of folk song or remodelled established folk ballads. Child included many minstrel ballads in his collection on the ground that fragments of traditional balladry were embedded in them. The blatant style of minstrelsy marks these ballads off sharply from folk creations. In violation of the strict impersonality of the folk ballads, minstrels constantly intrude into their narratives with moralizing comments and fervent assurances that they are not lying at the very moment when they are most fabulous. The

Difference
between
minstrel
and folk
ballads

minstrels manipulate the story with coarse explicitness, begging for attention in a servile way, predicting future events in the story and promising that it will be interesting and instructive, shifting scenes obtrusively, reflecting on the characters' motives with partisan prejudice. Often their elaborate performances are parcelled out in clear-cut divisions, usually called fits or cantos, in order to forestall tedium and build up suspense by delays and piecemeal revelations. Several of the surviving minstrel pieces are poems in praise of such noble houses as the Armstrongs ("Johnie Armstrong"), the Stanleys ("The Rose of England"), and the Percys ("The Battle of Otterburn," "The Hunting of the Cheviot," "The Earl of Westmoreland"), doubtless the work of propagandists in the employ of these families. The older Robin Hood ballads are also minstrel propaganda, glorifying the virtues of the yeomanry, the small independent landowners of preindustrial England. The longer, more elaborate minstrel ballads were patently meant to be recited rather than sung.

Broadside ballad. Among the earliest products of the printing press were broadsheets about the size of handbills on which were printed the text of ballads. A crude woodcut often headed the sheet, and under the title it was specified that the ballad was to be sung to the tune of some popular air. Musical notation seldom appeared on the broadsides; those who sold the ballads in the streets and at country fairs sang their wares so that anyone unfamiliar with the tune could learn it by listening a few times to the ballad-monger's rendition. From the 16th century until the end of the 19th century, broadsides, known also as street ballads, stall ballads, or slip songs, were a lively commodity, providing employment for a troop of hack poets. Before the advent of newspapers, the rhymed accounts of current events provided by the broadside ballads were the chief source of spectacular news. Every sensational public happening was immediately clapped into rhyme and sold on broadsheets. Few of the topical pieces long survived the events that gave them birth, but a good number of pathetic tragedies, such as "The Children in the Wood" and broadsides about Robin Hood, Guy of Warwick, and other national heroes, remained perennial favourites. Although the broadside ballad represents the adaptation of the folk ballad to the urban scene and middle class sensibilities, the general style more closely resembles minstrelsy, only with a generous admixture of vulgarized traits borrowed from book poetry. A few folk ballads appeared on broadsheets; many ballads, however, were originally broadside ballads the folk adapted.

Literary ballads. The earliest literary imitations of ballads were modelled on broadsides, rather than on folk ballads. In the early part of the 18th century, Jonathan Swift, who had written political broadsides in earnest, adapted the style for several jocular bagatelles. Poets such as Swift, Matthew Prior, and William Cowper in the 18th century and Thomas Hood, W.M. Thackeray, and Lewis Carroll in the 19th century made effective use of the jingling metres, forced rhymes, and unbuttoned style for humorous purposes. Lady Wardlaw's "Hardyknute" (1719), perhaps the earliest literary attempt at a folk ballad, was dishonestly passed off as a genuine product of tradition. After the publication of Thomas Percy's ballad compilation *Reliques of Ancient English Poetry* in 1765, ballad imitation enjoyed a considerable vogue, which properly belongs in the history of poetry rather than balladry.

SUBJECT MATTER

The supernatural. The finest of the ballads are deeply saturated in a mystical atmosphere imparted by the presence of magical appearances and apparatus. "The Wife of Usher's Well" laments the death of her children so unconsolably that they return to her from the dead as revenants; "Willie's Lady" cannot be delivered of her child because of her wicked mother-in-law's spells, an enchantment broken by a beneficent household spirit; "The Great Silkie of Sule Skerry" begets upon an "earthly" woman a son, who, on attaining maturity, joins his seal father in the sea, there shortly to be killed by his mother's human husband; "Kemp Owyne" disenchant a bespelled maiden by kissing her despite her bad breath and savage looks. An

encounter between a demon and a maiden occurs in "Lady Isabel and the Elf-Knight," the English counterpart of the ballads known to the Dutch-Flemish as "Herr Halewijn," to Germans as "Ulinger," to Scandinavians as "Kvindermorderen" and to the French as "Renaud le Tueur de Femme." In "The House Carpenter," a former lover (a demon in disguise) persuades a wife to forsake husband and children and come away with him, a fatal decision as it turns out. In American and in late British tradition the supernatural tends to get worked out of the ballads by being rationalized: instead of the ghost of his jilted sweetheart appearing to Sweet William of "Fair Margaret and Sweet William" as he lies in bed with his bride, it is rather the dead girl's image in a dream that kindles his fatal remorse. In addition to those ballads that turn on a supernatural occurrence, casual supernatural elements are found all through balladry.

Romantic tragedies. The separation of lovers through a misunderstanding or the opposition of relatives is perhaps the commonest ballad story. "Barbara Allen" is typical: Barbara cruelly spurns her lover because of an unintentional slight; he dies of lovesickness, she of remorse. The Freudian paradigm operates rigidly in ballads: fathers oppose the suitors of their daughters, mothers the sweethearts of their sons. Thus "The Douglas Tragedy"—the Danish "Ribold and Guldborg"—occurs when an eloping couple is overtaken by the girl's father and brothers or "Lady Maisry," pregnant by an English lord, is burned by her fanatically Scottish brother. Incest, frequent in ballads recorded before 1800 ("Lizie Wan," "The Bonny Hind"), is shunned by modern tradition.

Romantic comedies. The outcome of a ballad love affair is not always, though usually, tragic. But even when true love is eventually rewarded, such ballad heroines as "The Maid Freed from the Gallows" and "Fair Annie," among others, win through to happiness after such bitter trials that the price they pay seems too great. The course of romance runs hardly more smoothly in the many ballads, influenced by the cheap optimism of broadsides, where separated lovers meet without recognizing each other: the girl is told by the "stranger" of her lover's defection or death; her ensuing grief convinces him of her sincere love: he proves his identity and takes the joyful girl to wife. "The Bailiff's Daughter of Islington" is a classic of the type. Later tradition occasionally foists happy endings upon romantic tragedies: in the American "Douglas Tragedy" the lover is not slain but instead gets the irate father at his mercy and extorts a dowry from him. With marriage a consummation so eagerly sought in ballads, it is ironical that the bulk of humorous ballads deal with shrewish wives ("The Wife Wrapped in Wether's Skin") or gullible cuckolds ("Our Goodman").

Crime. Crime, and its punishment, is the theme of innumerable ballads: his sweetheart poisons "Lord Randal"; "Little Musgrave" is killed by Lord Barnard when he is discovered in bed with Lady Barnard, and the lady, too, is gorily dispatched. The murders of "Jim Fisk," Johnny of "Frankie and Johnny," and many other ballad victims are prompted by sexual jealousy. One particular variety of crime ballad, the "last goodnight," represents itself falsely to be the contrite speech of a criminal as he mounts the scaffold to be executed. A version of "Mary Hamilton" takes this form, which was a broadside device widely adopted by the folk. "Tom Dooley" and "Charles Guiteau," the scaffold confession of the assassin of Pres. James A. Garfield, are the best known American examples.

Medieval romance. Perhaps a dozen or so ballads derive from medieval romances. As in "Hind Horn" and "Thomas Rymer," only the climactic scene is excerpted for the ballad. In general, ballads from romances have not worn well in tradition because of their unpalatable fabulous elements, which the modern folk apparently regard as childish. Thus "Sir Lionel" becomes in America "Bangum and the Boar," a humorous piece to amuse children. Heterodox apocryphal legends that circulated widely in the Middle Ages are the source of almost all religious ballads, notable "Judas," "The Cherry-Tree Carol," and "The Bitter Withy." The distortion of biblical narrative is not peculiarly British: among others, the Russian bal-

Broadside ballad as a source of news

The Freudian paradigm

lads of Samson and Solomon, the Spanish "Pilgrim to Compostela" and the French and Catalonian ballads on the penance of Mary Magdalence reshape canonical stories radically.

Historical ballads. Historical ballads date mainly from the period 1550–750, though a few, like "The Battle of Otterburn," celebrate events of an earlier date, in this case 1388. "The Hunting of the Cheviot," recorded about the same time and dealing with the same campaign, is better known in a late broadside version called "Chevy Chase." The details in historical ballads are usually incorrect as to fact because of faulty memory or partisan alterations, but they are valuable in reflecting folk attitudes toward the events they imperfectly report. For example, neither "The Death of Queen Jane," about one of the wives of Henry VIII, nor "The Bonny Earl of Murray" is correct in key details, but they accurately express the popular mourning for these figures. By far the largest number of ballads that can be traced to historical occurrences have to do with local skirmishes and matters of regional rather than national importance. The troubled border between England and Scotland in the 16th and early 17th century furnished opportunities for intrepid displays of loyalty, courage, and cruelty that are chronicled in such dramatic ballads as "Edom o Gordon," "The Fire of Frenndraught," "Johnny Cock," "Johnnie Armstrong," and "Hobie Noble." Closely analogous to these are Spanish *romances* such as "The Seven Princes of Lara," on wars between Moors and Christians.

Disaster. Sensational shipwrecks, plagues, train wrecks, mine explosions—all kinds of shocking acts of God and man—were regularly chronicled in ballads, a few of which remained in tradition, probably because of some special charm in the language or the music. The shipwreck that lies in the background of one of the most poetic of all ballads "Sir Patrick Spens" cannot be fixed, but "The Titanic," "Casey Jones," "The Wreck on the C & O," and "The Johnstown Flood" are all circumstantially based on actual events.

Outlaws and badmen. Epic and saga heroes figure prominently in Continental balladries, notable examples being the Russian Vladimir, the Spanish Cid Campeador, the Greek Digenes Akritas, and the Danish Tord of Havsgaard and Diderik. This kind of hero never appears in English and Scottish ballads. But the outlaw hero of the type of the Serbian Makro Kraljević or the Danish Marsk Stig is exactly matched by the English Robin Hood, who is the hero of some 40 ballads, most of them of minstrel or broadside provenance. His chivalrous style and generosity to the poor was imitated by later ballad highwaymen in "Dick Turpin," "Brennan on the Moor," and "Jesse James." "Henry Martyn" and "Captain Kidd" were popular pirate ballads, but the most widely sung was "The Flying Cloud," a contrite "goodnight" warning young men to avoid the curse of piracy. The fact that so many folk heroes are sadistic bullies ("Stagolee"), robbers ("Dupree"), or pathological killers ("Sam Bass," "Billy the Kid") comments on the folk's hostile attitude toward the church, constabulary, banks, and railroads. The kindly, law-abiding, devout, enduring steel driver "John Henry" is a rarity among ballad heroes.

Occupational ballads. A large section of balladry, especially American, deals with the hazards of such occupations as seafaring ("The Greenland Whale Fishery"), lumbering ("The Jam on Gerry's Rock"), mining ("The Avondale Mine Disaster"), herding cattle ("Little Joc the Wrangler"), and the hardships of frontier life ("The Arkansaw Traveler"). But men in these occupations sang ballads also that had nothing to do with their proper work: "The Streets of Laredo," for example, is known in lumberjack and soldier versions as well as the usual cowboy lament version, and the pirate ballad "The Flying Cloud" was much more popular in lumbermen's shanties than in forecastles.

CHRONOLOGY

Singing stories in song, either stories composed for the occasion out of a repertory of traditional motifs or phrases or stories preserved by memory and handed down orally, is found in most primitive cultures. The ballad habit thus

is unquestionably very ancient. But the ballad genre itself could not have existed in anything like its present form before about 1100. "Judas," the oldest example found in Francis James Child's exhaustive collection, *The English and Scottish Popular Ballads* (1882–98), dates from 1300, but until the 17th century ballad records are sparse indeed. As an oral art, the ballad does not need to be written down to be performed or preserved; in any case, many of the carriers of the ballad tradition are illiterate and could not make use of a written and notated ballad. The few early ballads' records survived accidentally, due to some monk's, minstrel's, or antiquary's fascination with rustic pastimes.

The precise date of a ballad, therefore, or even any particular version of a ballad, is almost impossible to determine. In fact, to ask for the date of a folk ballad is to show that one misunderstands the peculiar nature of balladry. As remarked earlier, the first recording of a ballad must not be assumed to be the ballad's original form; behind each recorded ballad can be one detected the working of tradition upon some earlier form, since a ballad does not become a ballad until it has run a course in tradition. Historical ballads would seem on the surface to be easily datable, but their origins are usually quite uncertain. The ballad could have arisen long after the events it describes, basing itself, as do the Russian ballads of the Kievan cycle and the Spanish ballads about the Cid, on chronicles or popular legends. It is also likely that many historical ballads developed from the revamping of earlier ballads on similar themes through the alteration of names, places, and local details. (A.B.F.)

Precise dating almost impossible

Romance

Medieval romance developed in western Europe from the 12th century onward, and was extended up to and after the Renaissance, and reemerged in the 18th century.

The Old French word *romanz* originally meant "the speech of the people," or "the vulgar tongue," from a popular Latin word, *Romanice*, meaning written in the vernacular, in contrast with the written form of literary Latin. Its meaning then shifted from the language in which the work was written to the work itself. Thus, an adaptation of Geoffrey of Monmouth's *Historia regum Britanniae* (c. 1137), made by Wace of Jersey in 1155, was known as *Li Romanz de Bruu*, while an anonymous adaptation (of slightly later date) of Virgil's *Aeneid* was known as *Li Romanz d'Enéas*; it is difficult to tell whether in such cases *li romanz* still meant "the French version" or had already come to mean "the story." It soon specialized in the latter sense, however, and was applied to narrative compositions similar in character to those imitated from Latin sources but totally different in origin; and, as the nature of these compositions changed, the word itself acquired an increasingly wide spectrum of meanings. In modern French a *roman* is just a novel, whatever its content and structure; while in modern English the word "romance" (derived from Old French *romanz*) can mean either a medieval narrative composition or a love affair, or, again, a story about a love affair, generally one of a rather idyllic or idealized type, sometimes marked by strange or unexpected incidents and developments; and "to romance" has come to mean "to make up a story that has no connection with reality."

For a proper understanding of these changes it is essential to know something of the history of the literary form to which, since the Middle Ages, the term has been applied. The account that follows is intended to elucidate historically some of the ways in which the word is used in English and in other European languages.

THE COMPONENT ELEMENTS

The romances of love, chivalry, and adventure produced in 12th-century France have analogues elsewhere, notably in what are sometimes known as the Greek romances—narrative works in prose by Greek writers from the 1st century BC to the 3rd century AD. The first known, the fragmentary Ninus romance, in telling the story of the love of Ninus, mythical founder of Nineveh, anticipates the

Greek romances

Epic and
saga heroes

medieval *roman d'antiquité*. A number of works by writers of the 2nd and 3rd centuries AD—Chariton, Xenophon of Ephesus, Heliodorus, Achilles Tatius, and Longus—introduce a theme that was to reappear in the *roman d'aventure*: that of faithful lovers parted by accident or design and reunited only after numerous adventures. Direct connection, however, can be proved only in the case of the tale of *Apollonius of Tyre*, presumably deriving from a lost Greek original but known through a 3rd- or 4th-century Latin version. This too is a story of separation, adventure, and reunion, and, like the others (except for Longus' pastoral *Daphnis and Chloë*), it has a quasi-historical setting. It became one of the most popular and widespread stories in European literature during the Middle Ages and later provided Shakespeare with the theme of *Pericles*.

Style and subject matter. But the real debt of 12th-century romance to classical antiquity was incurred in a sphere outside that of subject matter. During the present century, scholars have laid ever-increasing emphasis on the impact of late classical antiquity upon the culture of medieval Europe, especially on that of medieval France. In particular, it is necessary to note the place that rhetoric (the systematic study of oratory) had assumed in the educational system of the late Roman Empire. Originally conceived as part of the training for public speaking, essential for the lawyer and politician, it had by this time become a literary exercise, the art of adorning or expanding a set theme: combined with grammar and enshrined in the educational system inherited by the Christian Church, rhetoric became an important factor in the birth of romance. Twelfth-century romance was, at the outset, the creation of "clerks"—professional writers who had been trained in grammar (that is to say, the study of the Latin language and the interpretation of Latin authors) and in rhetoric in the cathedral schools. They were skilled in the art of exposition, by which a subject matter was not only developed systematically but also given such meaning as the author thought appropriate. The "romance style" was, apparently, first used by the authors of three *romans d'antiquité*, all composed in the period 1150–65: *Roman de Thebes*, an adaptation of the epic *Thebais* by the late Latin poet Statius; *Roman d'Enéas*, adapted from Virgil's *Aeneid*; and *Roman de Troie*, a retelling by Benoît de Sainte-Maure of the tale of Troy, based not on Homer (who was not known in western Europe, where Greek was not normally read) but on 4th- and 5th-century Latin versions. In all three, style and subject matter are closely interconnected; elaborate set descriptions, in which the various features of what is described are gone through, item by item, and eulogized, result in the action's taking place in lavish surroundings, resplendent with gold, silver, marble, fine textiles, and precious stones. To these embellishments are added astonishing works of architecture and quaint technological marvels, that recall the Seven Wonders of the World and the reputed glories of Byzantium. *Troie* and *Enéas* have, moreover, a strong love interest, inspired by the Roman poet Ovid's conception of love as a restless malady. This concept produced the first portrayal in Western literature of the doubts, hesitations, and self-torture of young lovers, as exemplified in the Achilles-Polyxena story in *Troie* and in the Aeneas-Lavinia story in *Enéas*. Yet even more important is the way in which this new theme is introduced: the rhetorical devices appropriate to expounding an argument are here employed to allow a character in love to explore his own feelings, to describe his attitude to the loved one, and to explain whatever action he is about to take.

Developing psychological awareness. As W.P. Ker, a pioneer in the study of medieval epic and romance, observed in his *Epic and Romance* (1897), the advent of romance is "something as momentous and as far-reaching as that to which the name Renaissance is generally applied." The Old French poets who composed the chansons de geste (as the Old French epics are called) had been content to tell a story; they were concerned with statement, not with motivation, and their characters could act without explicitly justifying their actions. Thus, in what is one of the earliest and certainly the finest of the chansons de geste, the *Chanson de Roland* (c. 1100), the hero's decision to

fight on against odds—to let the rear guard of Charlemagne's army be destroyed by the Saracen hordes in the hopeless and heroic Battle of Roncevaux rather than sound his horn to call back Charlemagne—is not treated as a matter for discussion and analysis: the anonymous poet seems to take it for granted that the reader is not primarily concerned with the reason why things happened as they did. The new techniques of elucidating and elaborating material, developed by romance writers in the 12th century, produced a method whereby actions, motives, states of mind, were scrutinized and debated. The story of how Troilus fell in love with Briseïs and how, when taken to the Grecian camp, she deserted him for Diomedes (as related, and presumably invented, by Benoît de Sainte-Maure in his *Roman de Troie*) is not one of marvellous adventures in some exotic fairyland setting: it is clearly a theme of considerable psychological interest, and it was for this reason that it attracted three of the greatest writers of all time: Boccaccio in his *Filostrato* (c. 1338), Chaucer in his *Troilus and Criseyde* (before 1385), and Shakespeare in his *Troilus and Cressida* (c. 1601–02). With the 12th-century pioneers of what came to be called romance, the beginnings of the analytical method found in the modern novel can easily be recognized.

Sources and parallels. Where exactly medieval romance writers found their material when they were not simply copying classical or pseudo-classical models is still a highly controversial issue. Parallels to certain famous stories, such as that of Tristan and Iseult, have been found in regions as wide apart as Persia and Ireland: in the mid-11th-century Persian epic of *Wis and Ramin* and in the Old Irish *Diarmaid and Gráinne*; but while in the latter case it is possible to argue in favour of a genetic link between the two traditions, the former is more likely to be a case of parallel development due, on the one hand, to the inner logic of the theme and, on the other, to certain similarities in the ideological and social background of the two works. Failure to maintain the essential distinction between source and parallel has greatly hindered the understanding of the true nature of medieval romance and has led to the production of a vast critical literature the relevance of which to the study of the genre is at best questionable.

The marvellous. The marvellous is by no means an essential ingredient of "romance" in the sense in which it has been defined. Yet to most English readers the term romance does carry implications of the wonderful, the miraculous, the exaggerated, and the wholly ideal. Ker regarded much of the literature of the Middle Ages as "romantic" in this sense—the only types of narrative free from such "romanticizing" tendencies being the historical and family narrative, or Icelanders' sagas developed in classical Icelandic literature at the end of the 12th and in the early 13th century. The *Chanson de Roland* indulges freely in the fantastic and the unreal: hence Charlemagne's patriarchal age and preternatural strength (he is more than 200 years old when he conquers Spain); or the colossal numbers of those slain by the French; or, again, the monstrous races of men following the Saracen banners. Pious legends, saints' lives, and stories of such apocryphal adventures as those of the Irish St. Brendan (c. 486–578) who, as hero of a legend first written down in the 9th century, *Navigatio Brendani*, and later widely translated and adapted, wanders among strange islands on his way to the earthly paradise—these likewise favour the marvellous. The great 12th-century *Roman d'Alexandre*, a *roman d'antiquité* based on and developing the early Greek romance of Alexander the Great (the Alexander romance), was begun in the first years of the century by Alberic de Briançon and later continued by other poets. It introduces fantastic elements, more especially technological wonders and the marvels of India: the springs of rejuvenation, the flower-maidens growing in a forest, the cynocephali (dog-headed men), the bathyscaphe that takes Alexander to the bottom of the ocean, and the car in which he is drawn through the air by griffins on his celestial journey.

The setting. The fact that so many medieval romances are set in distant times and remote places is not an essential feature of romance but rather a reflection of its

The development of more sophisticated techniques

Earliest works in the "romance style"

Errors of
history and
geography
in early
romance

origins. As has been seen, the Old French word *romanz* early came to mean "historical work in the vernacular." All the *romans d'antiquité* have a historical or pseudo-historical theme, whether they evoke Greece, Troy or the legendary world of Alexander; but, while making some attempt to give antiquity an exotic aspect by means of marvels or technological wonders, medieval writers were quite unable to create a convincing historical setting; and thus in all important matters of social life and organization they projected the western European world of the 12th century back into the past. Similarly, historical and contemporary geography were not kept separate. The result is often a confused jumble, as, for example, in the Anglo-Norman Hue de Rotelande's *Protesilaus*, in which the characters have Greek names: the action takes place in Burgundy, Crete, Calabria, and Apulia; and Theseus is described as "king of Denmark." This lavish use of exotic personal and geographical names and a certain irresponsibility about settings was still to be found in some of Shakespeare's romantic comedies: the "seacoast of Bohemia" in *The Winter's Tale* is thoroughly medieval in its antecedents. In the medieval period, myth and folktale and straightforward fact were on an equal footing. Not that any marvel or preternatural happening taking place in secular (as opposed to biblical) history was necessarily to be believed: it was simply that the remote times and regions were convenient locations for picturesque and marvellous incidents. It is, indeed, at precisely this point that the transition begins from the concept of romance as "past history in the vernacular" to that of "a wholly fictitious story."

MEDIEVAL VERSE ROMANCES

Arthurian romance. *The matter of Britain.* In his *Historia regum Britanniae*, Geoffrey of Monmouth "invented history" by drawing on classical authors, the Bible, and Celtic tradition to create the story of a British kingdom, to some extent paralleling that of Israel. He described the rise of the British people to glory in the reigns of Uther Pendragon and Arthur, then the decline and final destruction of the kingdom, with the exile of the British survivors and their last king, Cadwalader. Romances that have Arthur or some of his knights as main characters were classified as *matière de Bretagne* by Jehan Bodel (fl. 1200) in a well-known poem. There is in this "matter of Britain" a certain amount of material ultimately based on the belief—probably Celtic in origin—in an otherworld into which men can penetrate, where they can challenge those who inhabit it or enjoy the love of fairy women. Such themes appear in a highly rationalized form in the lays (*lais*) of the late 12th-century Marie de France, although she mentions Arthur and his queen only in one, the lay of *Lanval*.

Chrétien de Troyes. It was Chrétien de Troyes (fl. 1165–80) who in five romances (*Erec*; *Cligès*; *Lancelot, ou Le Chevalier de la charrette*; *Yvain, ou Le Chevalier au lion*; and *Perceval, ou Le Conte du Graal*) fashioned a new type of narrative based on the matter of Britain. The internal debate and self-analysis of the *roman d'antiquité* is here used with artistry. At times, what seems to matter most to the poet is not the plot but the thematic pattern he imposes upon it and the significance he succeeds in conveying, either in individual scenes in which the action is interpreted by the characters in long monologues or through the work as a whole. In addition to this, he attempts what he himself calls a *conjointure*—that is, the organization into a coherent whole of a series of episodes. The adventures begin and end at the court of King Arthur; but the marvels that bring together material from a number of sources are not always meant to be believed, especially as they are somehow dovetailed into the normal incidents of life at a feudal court. Whatever Chrétien's intentions may have been, he inaugurated what may be called a Latin tradition of romance—clear, hard, bright, adorned with rhetoric, in which neither the courtly sentiment nor the enchantments are seriously meant. Chrétien had only one faithful follower, the trouvère Raoul de Houdenc (fl. 1200–30), author of *Méraugis de Portlesguez*. He shared Chrétien's taste for love casuistry, rhetorical adornment, and fantastic adventure. For both of these authors, elements of

rhetoric and self-analysis remain important, although the dose of rhetoric varies from one romance to another. Even in Chrétien's *Perceval, ou Le Conte du Graal* ("Perceval, or the Romance of the Grail")—the work in which the Grail appears for the first time in European literature—the stress is on narrative incident interspersed with predictions of future happenings and retrospective explanations. Arthurian romances of the period 1170–1250 are *romans d'aventure*, exploiting the strange, the supernatural, and the magical in the Arthurian tradition. A number (for example, *La Mule sans frein* ["The Mule Without a Bridle"], c. 1200, and *L'Âtre périlleux* ["The Perilous Churchyard"], c. 1250) have as their hero Arthur's nephew Gawain, who in the earlier Arthurian verse romances is a type of the ideal knight.

Love as a major theme. The treatment of love varies greatly from one romance to another. It is helpful to distinguish sharply here between two kinds of theme: the one, whether borrowed from classical antiquity (such as the story of Hero and Leander or that of Pyramus and Thisbe, taken from Ovid's *Metamorphoses*) or of much more recent origin, ending tragically; the other ending with marriage, reconciliation, or the reunion of separated lovers. It is noteworthy that "romance," as applied to a love affair in real life, has in modern English the connotation of a happy ending. This is also true of most Old French love romances in verse: the tragic ending is rare and is usually linked with the theme of the lover who, finding his or her partner dead, joins the beloved in death, either by suicide or from grief.

The Tristan story. The greatest tragic love story found as a romance theme is that of Tristan and Iseult. It was given the form in which it has become known to succeeding generations in about 1150–60 by an otherwise unknown Old French poet whose work, although lost, can be reconstructed in its essentials from surviving early versions based upon it. Probably closest in spirit to the original is the fragmentary version of c. 1170–90 by the Norman poet Bérout. From this it can be inferred that the archetypal poem told the story of an all-absorbing passion caused by a magic potion, a passion stronger than death yet unable to triumph over the feudal order to which the heroes belong. The story ended with Iseult's death in the embrace of her dying lover and with the symbol of two trees growing from the graves of the lovers and intertwining their branches so closely that they could never be separated. Most later versions, including a courtly version by an Anglo-Norman poet known only as Thomas, attempt to resolve the tragic conflict in favour of the sovereignty of passion and to turn the magic potion into a mere symbol. Gottfried von Strassburg's German version, *Tristan und Isolde* (c. 1210), based on Thomas, is one of the great courtly romances of the Middle Ages; but, although love is set up as the supreme value and as the object of the lovers' worship, the mellifluous and limpid verse translates the story into the idyllic mode. Another tragic and somewhat unreal story is that told in the anonymous *Chastelaine de Vergi* (c. 1250), one of the gems of medieval poetry, in which the heroine dies of grief because, under pressure, her lover has revealed their secret and adulterous love to the duke of Burgundy. The latter tells it to his own wife, who allows the heroine to think that her lover has betrayed her. The theme of the dead lover's heart served up by the jealous husband to the lady—tragic, sophisticated, and far-fetched—appears in the anonymous *Chastelain de Couci* (c. 1280) and again in *Das Herzmaere* by the late 13th-century German poet Konrad von Würzburg. The theme of the outwitting of the jealous husband, common in the fabliaux (short verse tales containing realistic, even coarse detail and written to amuse), is frequently found in 13th-century romance and in lighter lyric verse. It occurs both in the *Chastelain de Couci* and in the Provençal romance *Flamenca* (c. 1234), in which it is treated comically.

The theme of separation and reunion. But the theme that has left the deepest impress on romance is that of a happy resolution, after many trials and manifold dangers, of lovers' difficulties. As has been seen, this theme was derived from late classical Greek romance by way of *Apolonius of Tyre* and its numerous translations and variants.

Introduc-
tion of
the Grail
theme into
European
literature

Later
versions of
the Tristan
and Iseult
theme

A somewhat similar theme, used for pious edification, is that of the legendary St. Eustace, reputedly a high officer under the Roman emperor Trajan, who lost his position, property, and family only to regain them after many tribulations, trials, and dangers. The St. Eustace theme appears in *Guillaume d'Angleterre*, a pious tale rather than a romance proper, which some have attributed to Chrétien de Troyes.

The Floire and Blancheflor romance

A variant on the theme of separation and reunion is found in the romance of *Floire et Blancheflor* (c. 1170), in which Floire, son of the Saracen "king" of Spain, is parted by his parents from Blancheflor, daughter of a Christian slave of noble birth, who is sold to foreign slave dealers. He traces her to a tower where maidens destined for the sultan's harem are kept, and the two are reunited when he gains access to her there by hiding in a basket of flowers. This romance was translated into Middle High German, Middle Dutch, Norse, and Middle English (as *Floris and Blancheflor*, c. 1250) and in the early 13th century was imitated in *Aucassin et Nicolette*, which is a *chante-fable* (a story told in alternating sections of sung verse and recited prose) thought by some critics to share a common source with *Floire et Blancheflor*. In it, the roles and nationality, or religion, of the main characters are reversed; Nicolette, a Saracen slave converted to Christianity, who proves to be daughter of the king of Carthage, disguises herself as a minstrel in order to return to Aucassin, son of Count Gavin of Beaucaire. Jean Renart's *L'Escoufle* (c. 1200–02) uses the theme of lovers who, accidentally separated while fleeing together from the emperor's court, are eventually reunited; and the highly esteemed and influential *Guillaume de Palerne* (c. 1200) combines the theme of escaping lovers with that of the "grateful animal" (here a werewolf, which later resumes human shape as a king's son) assisting the lovers in their successful flight. The popular *Partenopeus de Blois* (c. 1180), of which 10 French manuscripts and many translated versions are known, resembles the Cupid and Psyche story told in the Roman writer Apuleius' *Golden Ass* (2nd century AD), although there is probably no direct connection. In the early 13th-century *Galeran de Bretagne*, Galeran loves Fresne, a founding brought up in a convent; the correspondence between the two is discovered, and Fresne is sent away but appears in Galeran's land just in time to prevent him from marrying her twin sister, Fleurie.

The theme of a knight who undertakes adventures to prove to his lady that he is worthy of her love is represented by a variety of romances including the *Ipomedon* (1174–90) of Hue de Rotelande and the anonymous mid-13th-century Anglo-Norman *Gui de Warewic*. Finally, there are many examples of the "persecuted heroine" theme; in one variety a person having knowledge of some "corporal sign"—a birthmark or mole—on a lady wagers with her husband that he will seduce her and offer proof that he has done so (this is sometimes called the "Imogen theme" from its use in Shakespeare's *Cymbeline*). The deceit is finally exposed and the lady's honour vindicated. In the early 13th-century *Guillaume de Dôle* by Jean Renart, the birthmark is a rose; and in the *Roman de Violette*, written after 1225 by Gerbert de Montreuil, it is a violet. Philippe de Beaumanoir's *La Manekine* (c. 1270), Jean Maillart's *La Contesse d'Anjou* (1361), and Chaucer's *Man of Law's Tale* (after 1387) all treat the theme of the tribulations of a wife falsely accused and banished but, after many adventures, reunited with her husband.

MEDIEVAL PROSE ROMANCES

Arthurian themes. The Arthurian prose romances arose out of the attempt, made first by Robert de Boron in the verse romances *Joseph d'Arimathie, ou le Roman de l'estoire dou Graal* and *Merlin* (c. 1190–1200), to combine the fictional history of the Holy Grail with the chronicle of the reign of King Arthur. Robert gave his story an allegorical meaning, related to the person and work of Christ. A severe condemnation of secular chivalry and courtly love characterize the Grail branch of the prose Lancelot-Grail, or Vulgate, cycle as well as some parts of the post-Vulgate "romance of the Grail" (after 1225); in the one case, Lancelot (here representing fallen human nature)

and, in the other, Balain (who strikes the Dolorous Stroke) are contrasted with Galahad, a type of the Redeemer. The conflict between earthly chivalry and the demands of religion is absent from the *Perlesvaus* (after 1230?), in which the hero Perlesvaus (that is, Perceval) has Christological overtones and in which the task of knighthood is to uphold and advance Christianity. A 13th-century prose *Tristan* (*Tristan de Léonois*), fundamentally an adaptation of the Tristan story to an Arthurian setting, complicates the love theme of the original with the theme of a love rivalry between Tristan and the converted Saracen Palamède and represents the action as a conflict between the treacherous villain King Mark and the "good" knight Tristan.

In the 14th century, when chivalry enjoyed a new vogue as a social ideal and the great orders of secular chivalry were founded, the romance writers, to judge from what is known of the voluminous *Perceforest* (written c. 1330 and still unpublished in its entirety), evolved an acceptable compromise between the knight's duty to his king, to his lady, and to God. Chivalry as an exalted ideal of conduct finds its highest expression in the anonymous Middle English *Sir Gawayne and the Grene Knight* (c. 1370), whose fantastic beheading scene (presumably taken from a lost French prose romance source) is made to illustrate the fidelity to the pledged word, the trust in God, and the unshakable courage that should characterize the knight.

Structure. The Vulgate *Lancelot-Grail* cycle displays a peculiar technique of interweaving that enables the author (or authors) to bring together a large number of originally independent themes. The story of Lancelot, of Arthur's kingdom, and the coming of Galahad (Lancelot's son) are all interconnected by the device of episodes that diverge, subdivide, join, and separate again, so that the work is a kind of interlocking whole, devoid of unity in the modern sense but forming as impregnable a structure as any revolving around a single centre. One of its most important features is its capacity for absorbing contrasting themes, such as the story of Lancelot's love for Guinevere, Arthur's queen, and the Quest of the Grail; another feature is its ability to grow through continuations or elaborations of earlier themes insufficiently developed. The great proliferation of prose romances at the end of the Middle Ages would have been impossible without this peculiarity of structure. Unlike any work that is wholly true to the Aristotelian principle of indivisibility and isolation (or organic unity), the prose romances satisfy the first condition, but not the second: internal cohesion goes with a tendency to seek connections with other similar compositions and to absorb an increasingly vast number of new themes. Thus the prose *Tristan* brings together the stories of Tristan and Iseult, the rise and fall of Arthur's kingdom, and the Grail Quest. It early gave rise to an offshoot, the romance of *Palamède* (before 1240), which deals with the older generation of Arthur's knights. A similar example of "extension backward" is the *Perceforest*, which associates the beginnings of knighthood in Britain with both Brutus the Trojan (reputedly Aeneas' grandson and the legendary founder of Britain) and Alexander the Great and makes its hero, Perceforest, live long before the Christian era.

LATER DEVELOPMENTS

The Arthurian prose romances were influential in both Italy and Spain; and this favoured the development in these countries of works best described as *romans d'aventure*, with their constantly growing interest in tournaments, enchantments, single combat between knights, love intrigues, and rambling adventures. In Italy, early prose compilations of Old French epic material from the Charlemagne cycle were subsequently assimilated to the other great bodies of medieval French narrative fiction and infused with the spirit of Arthurian prose romance. The great Italian heroic and romantic epics, Matteo Boiardo's *Orlando innamorato* (1483) and Ludovico Ariosto's *Orlando furioso* (1516), are based on this fusion. The serious themes of the Holy Grail and death of Arthur left no mark in Italy. The romantic idealism of Boiardo and Ariosto exploits instead the worldly adventures and the love sentiment of Arthurian prose romance, recounted lightly and with a sophisticated humour.

The ideals of chivalry

The special structure of the Lancelot-Grail romance

The "Imogen theme"

The Spanish romance

In Spain the significant development is the appearance, as early as the 14th, or even the 13th, century, of a native prose romance, the *Amadís de Gaula*. Arthurian in spirit but not in setting and with a freely invented episodic content, this work, in the form given to it by Garci Rodríguez de Montalvo in its first known edition of 1508, captured the imagination of the polite society of western Europe by its blend of heroic and incredible feats of arms and tender sentiment and by its exaltation of an idealized and refined concept of chivalry. Quickly translated and adapted into French, Italian, Dutch, and English and followed by numerous sequels and imitations in Spanish and Portuguese, it remained influential for more than four centuries, greatly affecting the outlook and sensibility of western society. Cervantes parodied the fashion inspired by *Amadís* in *Don Quixote* (1605); but his admiration for the work itself caused him to introduce many of its features into his own masterpiece, so that the spirit and the character of chivalric romance may be said to have entered into the first great modern novel.

Parallels with sculpture and painting

More important still for the development of the novel form was the use made by romance writers of the technique of multiple thematic structure and "interweaving" earlier mentioned. Like the great examples of Romanesque ornamental art, both sculptural and pictorial, the cyclic romances of the late Middle Ages, while showing a strong sense of cohesion, bear no trace whatever of the classical concept of subordination to a single theme: an excellent proof, if proof were needed, of the limited relevance of this concept in literary aesthetics. Even those romances which, like the *Amadís* and its ancestor, the French prose *Lancelot*, had one great figure as the centre of action, cannot be said to have progressed in any way toward the notion of the unity of theme.

The spread and popularity of romance literature. This is as true of medieval romances as of their descendants, including the French and the English 18th-century novel and the pastoral romance, which, at the time of the Renaissance, revived the classical traditions of pastoral poetry and led to the appearance, in 1504, of the *Arcadia* by the Italian poet Jacopo Sannazzaro and, in about 1559, of the *Diana* by the Spanish poet and novelist Jorge de Montemayor. Both works were widely influential in translation, and each has claims to be regarded as the first pastoral romance, but in spirit *Diana* is the true inheritor of the romance tradition, giving it, in alliance with the pastoral, a new impetus and direction.

The social milieu

Medieval romance began in the 12th century when clerks, working for aristocratic patrons, often ladies of royal birth such as Eleanor of Aquitaine and her daughters, Marie de Champagne and Matilda, wife of Henry the Lion, duke of Saxony, began to write for a leisured and refined society. Like the courtly lyric, romance was a vehicle of a new aristocratic culture which, based in France, spread to other parts of western Europe. Translations and adaptations of French romances appear early in German: the *Roman d'Enéas*, in a version written by Heinrich von Veldeke before 1186, and the archetypal Tristan romance in Eilhart von Oberg's *Tristan* of c. 1170–80. In England many French romances were adapted, sometimes very freely, into English verse and prose from the late 13th to the 15th century; but by far the most important English contribution to the development and popularization of romance was the adaptation of a number of French Arthurian romances completed by Sir Thomas Malory in 1469–70 and published in 1485 by William Caxton under the title of *Le Morte Darthur*. In the Scandinavian countries the connection with the Angevin rulers of England led to importation of French romances in the reign (1217–63) of Haakon of Norway.

The decline of romance. As has been seen, in the later Middle Ages the prose romances were influential in France, Italy, and Spain, as well as in England; and the advent of the printed book made them available to a still wider audience. But although they continued in vogue into the 16th century, with the spread of the ideals of the New Learning, the greater range and depth of vernacular literature, and the rise of the neoclassical critics, the essentially medieval image of the perfect knight was

bound to change into that of the scholar-courtier, who, as presented by the Italian Baldassare Castiglione in his *Il Cortegiano* (published 1528), embodies the highest moral ideals of the Renaissance. The new Spanish romances continued to enjoy international popularity until well into the 17th century and in France gave rise to compendious sentimental romances with an adventurous, pastoral, or pseudo-historical colouring popular with Parisian *salon* society until c. 1660. But the French intellectual climate, especially after the beginning of the so-called classical period in the 1660s, was unfavourable to the success of romance as a "noble" genre. Before disappearing, however, the romances lent the French form of their name to such romances as Antoine Furetière's *Le Roman bourgeois* (1666) and Paul Scarron's *Le Roman comique* (1651–57). These preserved something of the outward form of romance but little of its spirit; and while they transmitted the name to the kind of narrative fiction that succeeded them, they were in no sense intermediaries between its old and its new connotations. The great critical issue dominating the thought of western Europe from about 1660 onward was that of "truth" in literature; and romance, as being "unnatural" and unreasonable, was condemned. Only in England and Germany did it find a home with poets and novelists. Thus, while Robert Boyle, the natural philosopher, in his *Occasional Discourses* (1666) was inveighing against gentlemen whose libraries contained nothing more substantial than "romances," Milton, in *Paradise Lost*, could still invoke "what resounds/In fable or romance of Uther's son . . ."

The search for "truth" in literature

The 18th-century romantic revival. The 18th century in both England and Germany saw a strong reaction against the rationalistic canons of French classicism—a reaction that found its positive counterpart in such romantic material as had survived from medieval times. The Gothic romances, of which Horace Walpole's *Castle of Otranto* (1764; dated 1765) is the most famous, are perhaps of less importance than the ideas underlying the defense of romance by Richard Hurd in his *Letters on Chivalry and Romance* (1762). To Hurd, romance is not truth but a delightful and necessary holiday from common sense. This definition of romance (to which both Ariosto and Chrétien de Troyes would no doubt have subscribed) inspired on the one hand the romantic epic *Oberon* (1780) and on the other the historical romances of Sir Walter Scott. But influential though Scott's romantic novels may have been in every corner of Europe (including the Latin countries), it was the German and English Romantics who, with a richer theory of the imagination than Hurd's, were able to recapture something of the spirit and the structure of romance—the German Romantics by turning to their own medieval past; the English, by turning to the tradition perpetuated by Edmund Spenser and Shakespeare.

(E.Vt./F.Wh.)

Saga

In medieval Iceland the literary term saga denoted any kind of story or history in prose, irrespective of the kind or nature of the narrative or the purposes for which it was written. Used in this general sense, the term applies to a wide range of literary works, including those of hagiography (biography of saints), historiography, and secular fiction in a variety of modes. Lives of the saints and other stories for edification are entitled sagas, as are the Norse versions of French romances and the Icelandic adaptations of various Latin histories. Chronicles and other factual records of the history of Scandinavia and Iceland down to the 14th century are also included under the blanket term saga literature. In a stricter sense, however, the term saga is confined to legendary and historical fictions, in which the author has attempted an imaginative reconstruction of the past and organized the subject matter according to certain aesthetic principles. Using the distinctive features of the hero as principal guideline, medieval Icelandic narrative fiction can be classified as: (1) kings' sagas, (2) legendary sagas, and (3) sagas of Icelanders.

The origin and evolution of saga writing in Iceland are largely matters for speculation. A common pastime on

Icelandic farms, from the 12th century down to modern times, was the reading aloud of stories to entertain the household, known as *sagnaskemmtun* ("saga entertainment"). It seems to have replaced the traditional art of storytelling. All kinds of written narratives were used in *sagnaskemmtun*: secular, sacred, historical, and legendary. The Icelandic church took a sympathetic view of the writing and reading of sagas, and many of the authors whose identity is still known were monks or priests.

NONFICTIONAL SAGA LITERATURE

Translations. European narratives were known in Iceland in the 12th and 13th centuries and undoubtedly served as models for Icelandic writers when they set out to form a coherent picture of early Scandinavian history. Translations of lives of the saints and the apostles and accounts of the Holy Virgin testify to the skill of Icelandic prose writers in handling the vernacular for narrative purposes from the 12th century onward. Histories were also adapted and translated from Latin, based on those of the 7th- and 8th-century Anglo-Saxon writer Bede, the 7th-century Spanish historian Isidore of Seville, and others; on fictitious accounts of the Trojan wars, notably, one of the 5th century attributed to Dares Phrygius and one of the 4th century attributed to Dictys Cretensis; on the 12th-century British chronicler Geoffrey of Monmouth; and on the 1st-century Roman historians Sallust and Lucan. In the 13th century, Abbot Brandr Jónsson wrote a history of the Jews based on the Vulgate, on the 10th-century biblical scholar Peter Comestor, and on other sources.

In the 13th century, saga literature was also enriched by Norwegian prose translations of French romance literature. These soon found their way into Iceland, where they were popular and a strong influence on native storywriting. Probably the earliest, *Tristranis saga* (the story of Tristan and Iseult), was translated in 1226. Most of the themes of French romance appear in Icelandic versions; e.g., *Karlamagnús saga* was based on Charlemagne legends.

Native historical accounts. Icelandic historians seem to have started writing about their country's past toward the end of the 11th century. Saemundr Sigfússon, trained as a priest in France, wrote a Latin history of the kings of Norway, now lost but referred to by later authors. The first Icelandic to use the vernacular for historical accounts was Ari Thorgilsson, whose *Íslendingabók* (or *Libellus Islanðorum* [*The Book of the Icelanders*]) survives. It is a concise description of the course of Icelandic history from the beginning of the settlement (c. 870) to 1118. Ari seems to have written this book about 1125, but before that date he may already have compiled (in collaboration with Kolskeggr Asbjarnarson) the so-called *Landnámabók* ("Book of Settlements"), which lists the names and land claims of about 400 settlers. Because this work survives only in 13th- and 14th-century versions, it is impossible to tell how much of it is Ari's. Both books gave the Icelanders a clear picture of the beginning of their society; both works served to stimulate public interest in the period during which events recounted in the sagas of Icelanders (see below) are supposed to have taken place. Other factual accounts of the history of Iceland followed later: *Kristni saga* describes Iceland's conversion to Christianity about the end of the 10th century and the emergence of a national church. *Hungrvaka* ("The Appetizer") contains accounts of the lives of the first five bishops of Skálholt, from the mid-11th to the third quarter of the 12th century; the biographies of other prominent bishops are in the *Biskupa sögur*. Though some of these have a strong hagiographical flavour, others are soberly written and of great historical value. The period c. 1100–1264 is also dealt with in several secular histories, known collectively as *Sturlunga saga*, the most important of which is the *Íslendinga saga* ("The Icelanders' Saga") of Sturla Thórdarson, who describes in memorable detail the bitter personal and political feuds that marked the final episode in the history of the Icelandic commonwealth (c. 1200–64).

LEGENDARY AND HISTORICAL FICTION

Kings' sagas. After Saemundr Sigfússon, Icelandic and Norwegian authors continued to explore the history of

Scandinavia in terms of rulers and royal families, some of them writing in Latin and others in the vernacular. Broadly speaking, the kings' sagas fall into two distinct groups: contemporary (or near contemporary) biographies and histories of remoter periods. To the first group belonged a now-lost work, written in about 1170 by an Icelander called Eiríkr Oddsson, dealing with several 12th-century kings of Norway. *Sverris saga* describes the life of King Sverrir (reigned 1184–1202). The first part was written by Abbot Karl Jónsson under the supervision of the King himself, but it was completed (probably by the Abbot) in Iceland after Sverrir's death. Sturla Thórdarson wrote two royal biographies: *Haakonar saga* on King Haakon Haakonsson (c. 1204–63) and *Magnús saga* on his son and successor, Magnus VI Lawmender (Lagaboter; reigned 1263–80); of the latter only fragments survive. In writing these sagas Sturla used written documents as source material and, like Abbot Karl before him, he also relied on the accounts of eyewitnesses. Works on the history of the earlier kings of Norway include two Latin chronicles of Norwegian provenance, one of which was compiled c. 1180, and two vernacular histories, also written in Norway, the so-called *Ágrip* (c. 1190) and *Fagrskinna* (c. 1230). The Icelandic *Morkinskinna* (c. 1220) deals with the kings of Norway from 1047–1177; an outstanding feature of it is that it tells some brilliant stories of Icelandic poets and adventurers who visited the royal courts of Scandinavia.

The kings' sagas reached their zenith in the *Heimskringla*, or *Noregs konunga sögur* ("History of the Kings of Norway"), of Snorri Sturluson, which describes the history of the royal house of Norway from legendary times down to 1177. Snorri was a leading 13th-century Icelandic poet, who used as sources all the court poetry from the 9th century onward that was available to him. He also used many earlier histories of the kings of Norway and other written sources. *Heimskringla* is a supreme literary achievement that ranks Snorri Sturluson with the great writers of medieval Europe. He interpreted history in terms of personalities rather than politics, and many of his character portrayals are superbly drawn. Two of the early kings of Norway, Olaf Tryggvason (reigned 995–1000) and Olaf Haraldsson (Olaf the Saint; reigned 1015–30), received special attention from Icelandic antiquarians and authors. Only fragments of a 12th-century *Olafs saga helga* ("St. Olaf's Saga") survive; a 13th-century biography of the same king by Styrmir Kárason is also largely lost. (Snorri Sturluson wrote a brilliant saga of St. Olaf, rejecting some of the grosser hagiographical elements in his sources; this work forms the central part of his *Heimskringla*.) About 1190 a Benedictine monk, Oddr Snorrason, wrote a Latin life of Olaf Tryggvason, of which an Icelandic version still survives. A brother in the same monastery, Gunnlaugr Leifsson, expanded this biography, and his work was incorporated into later versions of *Olafs saga Tryggvasonar*. Closely related to the lives of the kings of Norway are *Fareyinga saga*, describing the resistance of Faeroese leaders to Norwegian interference during the first part of the 11th century, and *Orkneyinga saga*, dealing with the rulers of the earldom of Orkney from about 900 to the end of the 12th century. These two works were probably written about 1200. The history of the kings of Denmark from c. 940 to 1187 is told in *Knyflinga saga*.

Legendary sagas. The learned men of medieval Iceland took great pride in their pagan past and copied traditional poems on mythological and legendary themes. In due course some of these narrative poems served as the basis for sagas in prose. In his *Edda* (probably written c. 1225), Snorri Sturluson tells several memorable stories, based on ancient mythological poems, about the old gods of the North, including such masterpieces as the tragic death of Balder and the comic tale of Thor's journey to giantland. Snorri's book also contains a summary of the legendary Nibelungen cycle. (A much fuller treatment of the same theme is to be found in *Völsunga saga* and *Thidriks saga*, the latter composed in Norway and based on German sources.) Other Icelandic stories based on early poetic tradition include *Heidreks saga*; *Hrólf's saga kraka*, which has a certain affinity with the Old English poem *Beowulf*; *Hálfs saga ok Hálfsrekka*; *Gautreks saga*; and *Ásmundar*

Contemporary writings of remoter periods

Themes of French romance literature

saga kappabana, which tells the same story as the Old High German *Hildebrandslied*. The term legendary sagas also covers a number of stories the antecedents and models of which are not exclusively native. These sagas are set in what might be called the legendary heroic age at one level and also vaguely in the more recent Viking age at the other, the action taking place in Scandinavia and other parts of the Viking world, from Russia to Ireland, but occasionally also in the world of myth and fantasy. It is mostly through valour and heroic exploits that the typical hero's personality is realized. He is, however, often a composite character, for some of his features are borrowed from a later and more refined ethos than that of early Scandinavia. He is in fact the synthesis of Viking ideals on the one hand and of codes of courtly chivalry on the other. Of individual stories the following are notable: *Egils saga ok Asmundar*, which skillfully employs the flashback device; *Bósa saga ok Herrauds*, exceptional for its erotic elements; *Fridthjófs saga*, a romantic love story; *Hrólfs saga Gautrekssonar*; *Göngu-Hrólfs saga*; and *Halfadanar saga Eysteinnssonar*. There are many more. The legendary sagas are essentially romantic literature, offering an idealized picture of the remote past, and many of them are strongly influenced by French romance literature. In these sagas the main emphasis is on a lively narrative, entertainment being their primary aim and function. Some of the themes in the legendary sagas are also treated in the *Gesta Danorum* of the 12th-century Danish historian Saxo Grammaticus, who states that some of his informants for the legendary history of Denmark were Icelanders.

Sagas of Icelanders. In the late 12th century, Icelandic authors began to fictionalize the early part of their history (c. 900–1050), and a new literary genre was born: the sagas of Icelanders. Whereas the ethos of the kings' sagas and of the legendary sagas is aristocratic and their principal heroes warlike leaders, the sagas of Icelanders describe characters who are essentially farmers or farmers' sons or at least people who were socially not far above the author's public, and their conduct and motivation are measurable in terms of the author's own ethos. These authors constantly aimed at geographic, social, and cultural verisimilitude; they made it their business to depict life in Iceland as they had experienced it or as they imagined it had actually been in the past. Though a good deal of the subject matter was evidently derived from oral tradition and thus of historical value for the period described, some of the best sagas are largely fictional; their relevance to the authors' own times mattered perhaps no less than their incidental information about the past. An important aim of this literature was to encourage people to attain a better understanding of their social environment and a truer knowledge of themselves through studying the real and imagined fates of their forbears. A spirit of humanism, sometimes coloured by a fatalistic heroic outlook, pervades the narrative. The edificatory role, however, was never allowed to get out of hand or dominate the literary art; giving aesthetic pleasure remained the saga writer's primary aim and duty.

Nothing is known of the authorship of the sagas of Icelanders, and it has proved impossible to assign a definite date to many of them. It seems improbable that in their present form any of them could have been written before c. 1200. The period c. 1230–90 has been described as the golden age of saga writing because such masterpieces as *Egils saga*, *Víga-Glúms saga*, *Gísla saga*, *Eyrbyggja saga*, *Hrafnkels saga Freysgoda*, *Bandamanna saga*, *Hansa-Thóris saga*, and *Njáls saga* appear to have been written during that time. Although a number of sagas date from the 14th century, only one, *Grettis saga*, can be ranked with the classical ones.

The sagas of Icelanders can be subdivided into several categories according to the social and ethical status of the principal heroes. In some, the hero is a poet who sets out from the rural society of his native land in search of fame and adventure to become the retainer of the king of Norway or some other foreign ruler. Another feature of these stories is that the hero is also a lover. To this group belong some of the early-13th-century sagas, including *Kormáks saga*, *Hallfredar saga*, and *Bjarnar saga Hítuðelakappa*.

In *Gunnlaugs saga ormstungu*, which may have been written after the middle of the 13th century, the love theme is treated more romantically than in the others. *Fostbraeðra saga* ("The Blood-Brothers' Saga") describes two contrasting heroes: one a poet and lover, the other a ruthless killer. *Egils saga* offers a brilliant study of a complex personality—a ruthless Viking who is also a sensitive poet, a rebel against authority from early childhood who ends his life as a defenseless, blind old man. In several sagas the hero becomes an outlaw fighting a hopeless battle against the social forces that have rejected him. To this group belong *Hardar saga ok Hólmverja* and *Droplaugarsona saga*; but the greatest of the outlaw sagas are *Gísla saga*, describing a man who murders his own brother-in-law and whose sister reveals his dark secret; and *Grettis saga*, which deals with a hero of great talents and courage who is constantly fighting against heavy odds and is treacherously slain by an unscrupulous enemy.

Most of the sagas of Icelanders, however, are concerned with people who are fully integrated members of society, either as ordinary farmers or as farmers who also act as chieftains. *Hrafnkels saga* describes a chieftain who murders his shepherd, is then tortured and humiliated for his crime, and finally takes cruel revenge on one of his tormentors. The hero who gives his name to *Hansa-Thóris saga* is a man of humble background who makes money as a peddler and becomes a wealthy but unpopular landowner. His egotism creates trouble in the neighbourhood, and after he has set fire to one of the farmsteads, killing the farmer and the entire household, he is prosecuted and later put to death. *Ölkoфра thátr* (the term *thátr* is often used for a short story) and *Bandamanna saga* ("The Confederates' Saga") satirize chieftains who fail in their duty to guard the integrity of the law and try to turn other people's mistakes into profit for themselves. The central plot in *Laxdæla saga* is a love triangle, in which the jealous heroine forces her husband to kill his best friend. *Eyrbyggja saga* describes a complex series of feuds between several interrelated families; *Hávardar saga* is about an old farmer who takes revenge on his son's killer, the local chieftain; *Víga-Glúms saga* tells of a ruthless chieftain who commits several killings and swears an ambiguous oath in order to cover his guilt; while *Vansdæla saga* is the story of a noble chieftain whose last act is to help his killer escape.

In the sagas of Icelanders justice, rather than courage, is often the primary virtue, as might be expected in a literature that places the success of an individual below the welfare of society at large. This theme is an underlying one in *Njáls saga*, the greatest of all the sagas. It is a story of great complexity and richness, with a host of brilliantly executed character portrayals and a profound understanding of human strengths and weaknesses. Its structure is highly complex, but at its core is the tragedy of an influential farmer and sage who devotes his life to a hopeless struggle against the destructive forces of society but ends it inexorably when his enemies set fire to his house, killing his wife and sons with him. (He.P.)

Novel

The novel is a genre of fiction, and fiction may be defined as the art or craft of contriving, through the written word, representations of human life that instruct or divert or both. The various forms that fiction may take are best seen less as a number of separate categories than as a continuum or, more accurately, a cline, with some such brief form as the anecdote at one end of the scale and the longest conceivable novel at the other. When any piece of fiction is long enough to constitute a whole book, as opposed to a mere part of a book, then it may be said to have achieved novelhood. But this state admits of its own quantitative categories, so that a relatively brief novel may be termed a novella (or, if the insubstantiality of the content matches its brevity, a novelette), and a very long novel may overflow the banks of a single volume and become a *roman-fleuve*, or river novel. Length is very much one of the dimensions of the genre.

The term novel is a truncation of the Italian word *novella*

Aim of the sagas of Icelanders

Theme of *Njáls saga*

(from the plural of Latin *novellus*, a late variant of *novus*, meaning "new"), so that what is now, in most languages, a diminutive denotes historically the parent form. The *novella* was a kind of enlarged anecdote like those to be found in the 14th-century Italian classic Boccaccio's *Decameron*, each of which exemplifies the etymology well enough. The stories are little new things, novelties, freshly minted diversions, toys; they are not reworkings of known fables or myths, and they are lacking in weight and moral earnestness. It is to be noted that, despite the high example of novelists of the most profound seriousness, such as Tolstoy, Henry James, and Virginia Woolf, the term novel still, in some quarters, carries overtones of lightness and frivolity. And it is possible to decry a tendency to triviality in the form itself. The ode or symphony seems to possess an inner mechanism that protects it from aesthetic or moral corruption, but the novel can descend to shameful commercial depths of sentimentality or pornography. It is the purpose of this section to consider the novel not solely in terms of great art but also as an all-purpose medium catering for all the strata of literacy.

The novel
and the
epic
compared

Such early ancient Roman fiction as Petronius' *Satyricon* of the 1st century AD and Lucius Apuleius' *Golden Ass* of the 2nd century contain many of the popular elements that distinguish the novel from its nobler born relative the epic poem. In the fictional works, the medium is prose, the events described are unheroic, the settings are streets and taverns, not battlefields and palaces. There is more low fornication than princely combat; the gods do not move the action; the dialogue is homely rather than aristocratic. It was, in fact, out of the need to find—in the period of Roman decline—a literary form that was anti-epic in both substance and language that the first prose fiction of Europe seems to have been conceived. The most memorable character in Petronius is a *nouveau riche* vulgarian; the hero of Lucius Apuleius is turned into a donkey; nothing less epic can well be imagined.

The medieval chivalric romance (from a popular Latin word, probably *Romanice*, meaning written in the vernacular, not in traditional Latin) restored a kind of epic view of man—though now as heroic Christian, not heroic pagan. At the same time, it bequeathed its name to the later genre of continental literature, the novel, which is known in French as *roman*, in Italian as *romanzo*, etc. (The English term romance, however, carries a pejorative connotation.) But that later genre achieved its first great flowering in Spain at the beginning of the 17th century in an antichivalric comic masterpiece—the *Don Quixote* of Cervantes, which, on a larger scale than the *Satyricon* or *The Golden Ass*, contains many of the elements that have been expected from prose fiction ever since. Novels have heroes, but not in any classical or medieval sense. As for the novelist, he must, in the words of the contemporary British-American W.H. Auden,

Become the whole of boredom, subject to
Vulgar complaints like love, among the Just
Be just, among the Filthy filthy too,
And in his own weak person, if he can,
Must suffer dully all the wrongs of Man.

The novel attempts to assume those burdens of life that have no place in the epic poem and to see man as unheroic, unredeemed, imperfect, even absurd. This is why there is room among its practitioners for writers of hard-boiled detective thrillers such as the contemporary American Mickey Spillane or of sentimental melodramas such as the prolific 19th-century English novelist Mrs. Henry Wood, but not for one of the unremitting elevation of outlook of a John Milton.

The reader may also be interested in the analogous treatment of a comparable genre in the section *Short story* which follows. Critical approaches are discussed in the section *Literary criticism* below.

ELEMENTS

Plot. The novel is propelled through its hundred or thousand pages by a device known as the story or plot. This is frequently conceived by the novelist in very simple terms, a mere nucleus, a jotting on an old envelope: for example, Charles Dickens' *Christmas Carol* (1843)

might have been conceived as "a misanthrope is reformed through certain magical visitations on Christmas Eve," or Jane Austen's *Pride and Prejudice* (1813) as "a young couple destined to be married have first to overcome the barriers of pride and prejudice," or Fyodor Dostoyevsky's *Crime and Punishment* (1866) as "a young man commits a crime and is slowly pursued in the direction of his punishment." The detailed working out of the nuclear idea requires much ingenuity, since the plot of one novel is expected to be somewhat different from that of another, and there are very few basic human situations for the novelist to draw upon. The dramatist may take his plot ready-made from fiction or biography—a form of theft sanctioned by Shakespeare—but the novelist has to produce what look like novelties.

Original
and
borrowed
plots

The example of Shakespeare is a reminder that the ability to create an interesting plot, or even any plot at all, is not a prerequisite of the imaginative writer's craft. At the lowest level of fiction, plot need be no more than a string of stock devices for arousing stock responses of concern and excitement in the reader. The reader's interest may be captured at the outset by the promise of conflicts or mysteries or frustrations that will eventually be resolved, and he will gladly—so strong is his desire to be moved or entertained—suspend criticism of even the most trite modes of resolution. In the least sophisticated fiction, the knots to be untied are stringently physical, and the denouement often comes in a sort of triumphant violence. Serious fiction prefers its plots to be based on psychological situations, and its climaxes come in new states of awareness—chiefly self-knowledge—on the parts of the major characters.

Melodramatic plots, plots dependent on coincidence or improbability, are sometimes found in even the most elevated fiction; E.M. Forster's *Howards End* (1910) is an example of a classic British novel with such a plot. But the novelist is always faced with the problem of whether it is more important to represent the formlessness of real life (in which there are no beginnings and no ends and very few simple motives for action) or to construct an artifact as well balanced and economical as a table or chair; since he is an artist, the claims of art, or artifice, frequently prevail.

There are, however, ways of constructing novels in which plot may play a desultory part or no part at all. The traditional picaresque novel—a novel with a rogue as its central character—like Alain Lesage's *Gil Blas* (1715) or Henry Fielding's *Tom Jones* (1749), depends for movement on a succession of chance incidents. In the works of Virginia Woolf, the consciousness of the characters, bounded by some poetic or symbolic device, sometimes provides all the fictional material. Marcel Proust's great *roman-fleuve*, *À la recherche du temps perdu* (1913–27; *Remembrance of Things Past*), has a metaphysical framework derived from the time theories of the philosopher Henri Bergson, and it moves toward a moment of truth that is intended to be literally a revelation of the nature of reality. Strictly, any scheme will do to hold a novel together—raw action, the hidden syllogism of the mystery story, prolonged solipsist contemplation—so long as the actualities or potentialities of human life are credibly expressed, with a consequent sense of illumination, or some lesser mode of artistic satisfaction, on the part of the reader.

Character. The inferior novelist tends to be preoccupied with plot; to the superior novelist the convolutions of the human personality, under the stress of artfully selected experience, are the chief fascination. Without character it was once accepted that there could be no fiction. In the period since World War II, the creators of what has come to be called the French *nouveau roman* (i.e., new novel) have deliberately demoted the human element, claiming the right of objects and processes to the writer's and reader's prior attention. Thus, in books termed *chosiste* (literally "thing-ist"), they make the furniture of a room more important than its human incumbents. This may be seen as a transitory protest against the long predominance of character in the novel, but, even on the popular level, there have been indications that readers can be held by things as much as by characters. Henry James could be

Objects as
"characters"

vague in *The Ambassadors* (1903) about the provenance of his chief character's wealth; if he wrote today he would have to give his readers a tour around the factory or estate. The popularity of much undistinguished but popular fiction has nothing to do with its wooden characters; it is machines, procedures, organizations that draw the reader. The success of Ian Fleming's British spy stories in the 1960s had much to do with their hero, James Bond's car, gun, and preferred way of mixing a martini.

But the true novelists remain creators of characters—prehuman, such as those in William Golding's *Inheritors* (1955); animal, as in Henry Williamson's *Tarka the Otter* (1927) or Jack London's *Call of the Wild* (1903); caricatures, as in much of Dickens; or complex and unpredictable entities, as in Tolstoy, Dostoyevsky, or Henry James. The reader may be prepared to tolerate the most wanton-seeming stylistic tricks and formal difficulties because of the intense interest of the central characters in novels as diverse as James Joyce's *Ulysses* (1922) and *Finnegans Wake* (1939) and Laurence Sterne's *Tristram Shandy* (1760–67).

Characters
as symbols

It is the task of literary critics to create a value hierarchy of fictional character, placing the complexity of the Shakespearean view of man—as found in the novels of Tolstoy and Joseph Conrad—above creations that may be no more than simple personifications of some single characteristic, like some of those by Dickens. It frequently happens, however, that the common reader prefers surface simplicity—easily memorable cartoon figures like Dickens' never-despairing Mr. Micawber and devious Uriah Heep—to that wider view of personality, in which character seems to engulf the reader, subscribed to by the great novelists of France and Russia. The whole nature of human identity remains in doubt, and writers who voice that doubt—like the French exponents of the *nouveau roman* Alain Robbe-Grillet and Nathalie Sarraute, as well as many others—are in effect rejecting a purely romantic view of character. This view imposed the author's image of himself—the only human image he properly possessed—on the rest of the human world. For the unsophisticated reader of fiction, any created personage with a firm position in time-space and the most superficial parcel of behavioral (or even sartorial) attributes will be taken for a character. Though the critics may regard it as heretical, this tendency to accept a character is in conformity with the usages of real life. The average person has at least a suspicion of his own complexity and inconsistency of makeup, but he sees the rest of the world as composed of much simpler entities. The result is that novels whose characters are created out of the author's own introspection are frequently rejected as not "true to life." But both the higher and the lower orders of novel readers might agree in condemning a lack of memorability in the personages of a work of fiction, a failure on the part of the author to seem to add to the reader's stock of remembered friends and acquaintances. Characters that seem, on recollection, to have a life outside the bounds of the books that contain them are usually the ones that earn their creators the most regard. Depth of psychological penetration, the ability to make a character real as oneself, seems to be no primary criterion of fictional talent.

Scene, or setting. The makeup and behaviour of fictional characters depend on their environment quite as much as on the personal dynamic with which their author endows them: indeed, in Émile Zola, environment is of overriding importance, since he believed it determined character. The entire action of a novel is frequently determined by the locale in which it is set. Thus, Gustave Flaubert's *Madame Bovary* (1857) could hardly have been placed in Paris, because the tragic life and death of the heroine have a great deal to do with the circumscriptions of her provincial milieu. But it sometimes happens that the main locale of a novel assumes an importance in the reader's imagination comparable to that of the characters and yet somehow separable from them. Wessex is a giant brooding presence in Thomas Hardy's novels, whose human characters would probably not behave much differently if they were set in some other rural locality of England. The popularity of Sir Walter Scott's "Waverley"

novels is due in part to their evocation of a romantic Scotland. Setting may be the prime consideration of some readers, who can be drawn to Conrad because he depicts life at sea or in the East Indies; they may be less interested in the complexity of human relationships that he presents.

Regional
novels

The regional novel is a recognized species. The sequence of four novels that Hugh Walpole began with *Rogue Herries* (1930) was the result of his desire to do homage to the part of Cumberland, in England, where he had elected to live. The great Yoknapatawpha cycle of William Faulkner, a classic of 20th-century American literature set in an imaginary county in Mississippi, belongs to the category as much as the once-popular confections about Sussex that were written about the same time by the English novelist Sheila Kaye-Smith. Many novelists, however, gain a creative impetus from avoiding the same setting in book after book and deliberately seeking new locales. The English novelist Graham Greene apparently needed to visit a fresh scene in order to write a fresh novel. His ability to encapsulate the essence of an exotic setting in a single book is exemplified in *The Heart of the Matter* (1948); his contemporary Evelyn Waugh stated that the West Africa of that book replaced the true remembered West Africa of his own experience. Such power is not uncommon: the Yorkshire moors have been romanticized because Emily Brontë wrote of them in *Wuthering Heights* (1847), and literary tourists have visited Stoke-on-Trent, in northern England, because it comprises the "Five Towns" of Arnold Bennett's novels of the early 20th century. Others go to the Monterey, California, of John Steinbeck's novels in the expectation of experiencing a *frisson* added to the locality by an act of creative imagination. James Joyce, who remained inexhaustibly stimulated by Dublin, has exalted that city in a manner that even the guidebooks recognize.

The setting of a novel is not always drawn from a real-life locale. The literary artist sometimes prides himself on his ability to create the totality of his fiction—the setting as well as the characters and their actions. In the Russian expatriate Vladimir Nabokov's *Invitation of a Friend* (1929) there is an entirely new space-time continuum, and the English scholar J.R.R. Tolkien in his *The Lord of the Rings* (1954–55) created an "alternative world" that appeals greatly to many who are dissatisfied with the existing one. The world of interplanetary travel was imaginatively created long before the first moon landing. The properties of the future envisaged by H.G. Wells's novels or by Aldous Huxley in *Brave New World* (1932) are still recognized in an age that those authors did not live to see. The composition of place can be a magical fictional gift.

Whatever the locale of his work, every true novelist is concerned with making a credible environment for his characters, and this really means a close attention to sense data—the immediacies of food and drink and colour—far more than abstractions like "nature" and "city." The London of Charles Dickens is as much incarnated in the smell of wood in lawyers' chambers as in the skyline and vistas of streets.

Narrative method and point of view. Where there is a story, there is a storyteller. Traditionally, the narrator of the epic and mock-epic alike acted as an intermediary between the characters and the reader; the method of Fielding is not very different from the method of Homer. Sometimes the narrator boldly imposed his own attitudes; always he assumed an omniscience that tended to reduce the characters to puppets and the action to a predetermined course with an end implicit in the beginning. Many novelists have been unhappy about a narrative method that seems to limit the free will of the characters, and innovations in fictional technique have mostly sought the objectivity of the drama, in which the characters appear to work out their own destinies without prompting from the author.

The epistolary method, most notably used by Samuel Richardson in *Pamela* (1740) and by Jean-Jacques Rousseau in *La nouvelle Héloïse* (1761), has the advantage of allowing the characters to tell the story in their own words, but it is hard to resist the uneasy feeling that a kind of divine editor is sorting and ordering the letters into his own pattern. The device of making the

narrator also a character in the story has the disadvantage of limiting the material available for the narration, since the narrator-character can know only those events in which he participates. There can, of course, be a number of secondary narratives enclosed in the main narrative, and this device—though it sometimes looks artificial—has been used triumphantly by Conrad and, on a lesser scale, by W. Somerset Maugham. A, the main narrator, tells what he knows directly of the story and introduces what B and C and D have told him about the parts that he does not know.

Seeking the most objective narrative method of all, Ford Madox Ford used, in *The Good Soldier* (1915), the device of the storyteller who does not understand the story he is telling. This is the technique of the “unreliable observer.” The reader, understanding better than the narrator, has the illusion of receiving the story directly. Joyce, in both his major novels, uses different narrators for the various chapters. Most of them are unreliable, and some of them approach the impersonality of a sort of disembodied parody. In *Ulysses*, for example, an episode set in a maternity hospital is told through the medium of a parodic history of English prose style. But, more often than not, the sheer ingenuity of Joyce’s techniques draws attention to the manipulator in the shadows. The reader is aware of the author’s cleverness where he should be aware only of the characters and their actions. The author is least noticeable when he is employing the stream of consciousness device, by which the inchoate thoughts and feelings of a character are presented in interior monologue—apparently unedited and sometimes deliberately near-unintelligible. It is because this technique seems to draw fiction into the psychoanalyst’s consulting room (presenting the raw material of either art or science, but certainly not art itself), however, that Joyce felt impelled to impose the shaping devices referred to above. Joyce, more than any novelist, sought total objectivity of narration technique but ended as the most subjective and idiosyncratic of stylists.

The problem of a satisfactory narrative point of view is, in fact, nearly insoluble. The careful exclusion of comment, the limitation of vocabulary to a sort of reader’s lowest common denominator, the paring of style to the absolute minimum—these puritanical devices work well for an Ernest Hemingway (who, like Joyce, remains, nevertheless, a highly idiosyncratic stylist) but not for a novelist who believes that, like poetry, his art should be able to draw on the richness of word play, allusion, and symbol. For even the most experienced novelist, each new work represents a struggle with the unconquerable task of reconciling all-inclusion with self-exclusion. It is noteworthy that Cervantes, in *Don Quixote*, and Nabokov, in *Lolita* (1955), join hands across four centuries in finding most satisfactory the device of the fictitious editor who presents a manuscript story for which he disclaims responsibility. But this highly useful method presupposes in the true author a scholarly, or pedantic, faculty not usually associated with novelists.

Scope, or dimension. No novel can theoretically be too long, but if it is too short it ceases to be a novel. It may or may not be accidental that the novels most highly regarded by the world are of considerable length—Cervantes’ *Don Quixote*, Dostoyevsky’s *Brothers Karamazov*, Tolstoy’s *War and Peace*, Dickens’ *David Copperfield*, Proust’s *À la recherche du temps perdu*, and so on. On the other hand, since World War II, brevity has been regarded as a virtue in works like the later novels of the Irish absurdist author Samuel Beckett and the *ficciones* of the Argentine Jorge Luis Borges, and it is only an aesthetic based on bulk that would diminish the achievement of Ronald Firbank’s short novels of the post-World War I era or the Evelyn Waugh who wrote *The Loved One* (1948). It would seem that there are two ways of presenting human character—one, the brief way, through a significant episode in the life of a personage or group of personages; the other, which admits of limitless length, through the presentation of a large section of a life or lives, sometimes beginning with birth and ending in old age. The plays of Shakespeare show that a full delineation of character can be effected in a very brief compass, so that, for this aspect of the novel,

length confers no special advantage. Length, however, is essential when the novelist attempts to present something bigger than character—when, in fact, he aims at the representation of a whole society or period of history.

No other cognate art form—neither the epic poem nor the drama nor the film—can match the resources of the novel when the artistic task is to bring to immediate, sensuous, passionate life the somewhat impersonal materials of the historian. *War and Peace* is the great triumphant example of the panoramic study of a whole society—that of early 19th-century Russia—which enlightens as the historian enlightens and yet also conveys directly the sensations and emotions of living through a period of cataclysmic change. In the 20th century, another Russian, Boris Pasternak, in his *Doctor Zhivago* (1957), expressed—though on a less than Tolstoyan scale—the personal immediacies of life during the Russian Revolution. Though of much less literary distinction than either of these two books, Margaret Mitchell’s *Gone with the Wind* (1936) showed how the American Civil War could assume the distanced pathos, horror, and grandeur of any of the classic struggles of the Old World.

Needless to say, length and weighty subject matter are no guarantee in themselves of fictional greatness. Among American writers, for example, James Jones’s celebration of the U.S. Army on the eve of World War II in *From Here to Eternity* (1951), though a very ambitious project, repels through indifferent writing and sentimental characterization; Norman Mailer’s *Naked and the Dead* (1948), an equally ambitious military novel, succeeds much more because of a tautness, a concern with compression, and an astringent objectivity that Jones was unable to match. Frequently the size of a novel is too great for its subject matter—as with Marguerite Young’s *Miss MacIntosh, My Darling* (1965), reputedly the longest single-volume novel of the 20th century, John Barth’s *Giles Goat-Boy* (1966), and John Fowles’s *Magus* (1965). Diffuseness is the great danger in the long novel, and diffuseness can mean slack writing, emotional self-indulgence, sentimentality.

Even the long picaresque novel—which, in the hands of a Fielding or his contemporary Tobias Smollett, can rarely be accused of sentimentality—easily betrays itself into such acts of self-indulgence as the multiplication of incident for its own sake, the coy digression, the easygoing jogtrot pace that subdues the sense of urgency that should lie in all fiction. If Tolstoy’s *War and Peace* is a greater novel than Fielding’s *Tom Jones* or Dickens’ *David Copperfield*, it is not because its theme is nobler, or more pathetic, or more significant historically; it is because Tolstoy brings to his panoramic drama the compression and urgency usually regarded as the monopolies of briefer fiction.

Sometimes the scope of a fictional concept demands a technical approach analogous to that of the symphony in music—the creation of a work in separate books, like symphonic movements, each of which is intelligible alone but whose greater intelligibility depends on the theme and characters that unify them. The French author Romain Rolland’s *Jean-Christophe* (1904–12) sequence is, very appropriately since the hero is a musical composer, a work in four movements. Among works of English literature, Lawrence Durrell’s *Alexandria Quartet* (1957–60) insists in its very title that it is a tetralogy rather than a single large entity divided into four volumes; the concept is “relativist” and attempts to look at the same events and characters from four different viewpoints. Anthony Powell’s *Dance to the Music of Time*, a multivolume series of novels that began in 1951 (collected 1962), may be seen as a study of a segment of British society in which the chronological approach is eschewed, and events are brought together in one volume or another because of a kind of parachronic homogeneity. C.P. Snow’s *Strangers and Brothers*, a comparable series that began in 1940 and continued to appear throughout the ’50s and into the ’60s, shows how a fictional concept can be realized only in the act of writing, since the publication of the earlier volumes antedates the historical events portrayed in later ones. In other words, the author could not know what the subject matter of the sequence would be until he was in sight of its end. Behind all these works lies the giant example of

The presence of the author

Works in several books

Proust's *roman-fleuve*, whose length and scope were properly coterminous with the author's own life and emergent understanding of its pattern.

Myth, symbolism, significance. The novelist's conscious day-to-day preoccupation is the setting down of incident, the delineation of personality, the regulation of exposition, climax, and denouement. The aesthetic value of the work is frequently determined by subliminal forces that seem to operate independently of the writer, investing the properties of the surface story with a deeper significance. A novel will then come close to myth, its characters turning into symbols of permanent human states or impulses, particular incarnations of general truths perhaps only realized for the first time in the act of reading. The ability to perform a quixotic act antedated *Don Quixote*, just as *bovarysme* existed before Flaubert found a name for it.

But the desire to give a work of fiction a significance beyond that of the mere story is frequently conscious and deliberate, indeed sometimes the primary aim. When a novel—like Joyce's *Ulysses* or John Updike's *Centaur* (1963) or Anthony Burgess' *Vision of Battlements* (1965)—is based on an existing classical myth, there is an intention of either ennobling a lowly subject matter, satirizing a debased set of values by referring them to a heroic age, or merely providing a basic structure to hold down a complex and, as it were, centrifugal picture of real life. Of *Ulysses* Joyce said that his Homeric parallel (which is worked out in great and subtle detail) was a bridge across which to march his 18 episodes; after the march the bridge could be "blown skyhigh." But there is no doubt that, through the classical parallel, the account of an ordinary summer day in Dublin is given a richness, irony, and universality unattainable by any other means.

The mythic or symbolic intention of a novel may manifest itself less in structure than in details which, though they appear naturalistic, are really something more. The shattering of the eponymous golden bowl in Henry James's 1904 novel makes palpable, and hence truly symbolic, the collapse of a relationship. Even the choice of a character's name may be symbolic. Sammy Mountjoy, in William Golding's *Free Fall* (1959), has fallen from the grace of heaven, the mount of joy, by an act of volition that the title makes clear. The eponym of *Doctor Zhivago* is so called because his name, meaning "The Living," carries powerful religious overtones. In the Russian version of the Gospel According to St. Luke, the angels ask the women who come to Christ's tomb: "*Chto vy ischlyote zhivago mezhdu myortvykh?*"—"Why do you seek the living among the dead?" And his first name, Yuri, the Russian equivalent of George, has dragon-slaying connotations.

The symbol, the special significance at a subnarrative level, works best when it can fit without obtrusion into a context of naturalism. The optician's trade sign of a huge pair of spectacles in F. Scott Fitzgerald's *Great Gatsby* (1925) is acceptable as a piece of scenic detail, but an extra dimension is added to the tragedy of Gatsby, which is the tragedy of a whole epoch in American life, when it is taken also as a symbol of divine myopia. Similarly, a cinema poster in Malcolm Lowry's *Under the Volcano* (1947), advertising a horror film, can be read as naturalistic background, but it is evident that the author expects the illustrated fiend—a concert pianist whose grafted hands are those of a murderer—to be seen also as a symbol of Nazi infamy: the novel is set at the beginning of World War II, and the last desperate day of the hero, Geoffrey Firmin, stands also for the collapse of Western civilization.

There are symbolic novels whose infranarrative meaning cannot easily be stated, since it appears to subsist on an unconscious level. Herman Melville's *Moby Dick* (1851) is such a work, as is D.H. Lawrence's novella *St. Mawr* (1925), in which the significance of the horse is powerful and mysterious.

USES

Interpretation of life. Novels are not expected to be didactic, like tracts or morality plays; nevertheless, in varying degrees of implicitness, even the "purest" works of fictional art convey a philosophy of life. The novels of Jane Austen, designed primarily as superior entertainment,

imply a desirable ordered existence, in which the comfortable decorum of an English rural family is disturbed only by a not-too-serious shortage of money, by love affairs that go temporarily wrong, and by the intrusion of self-centred stupidity. The good, if unrewarded for their goodness, suffer from no permanent injustice. Life is seen, not only in Jane Austen's novels but in the whole current of bourgeois Anglo-American fiction, as fundamentally reasonable and decent. When wrong is committed, it is usually punished, thus fulfilling Miss Prism's summation in Oscar Wilde's play *The Importance of Being Earnest* (1895), to the effect that in a novel the good characters end up happily and the bad characters unhappily: "that is why it is called fiction."

That kind of fiction called realistic, which has its origins in 19th-century France, chose the other side of the coin, showing that there was no justice in life and that the evil and the stupid must prevail. In the novels of Thomas Hardy there is a pessimism that may be taken as a corrective of bourgeois Panglossianism—the philosophy that everything happens for the best, satirized in Voltaire's *Candide* (1759)—since the universe is presented as almost impossibly malevolent. This tradition is regarded as morbid, and it has been deliberately ignored by most popular novelists. The "Catholic" novelists—such as François Mauriac in France, Graham Greene in England, and others—see life as mysterious, full of wrong and evil and injustice inexplicable by human canons but necessarily acceptable in terms of the plans of an inscrutable God. Between the period of realistic pessimism, which had much to do with the agnosticism and determinism of 19th-century science, and the introduction of theological evil into the novel, writers such as H.G. Wells attempted to create a fiction based on optimistic liberalism. As a reaction, there was the depiction of "natural man" in the novels of D.H. Lawrence and Ernest Hemingway.

For the most part, the view of life common to American and European fiction since World War II posits the existence of evil—whether theological or of that brand discovered by the French Existentialists, particularly Jean-Paul Sartre—and assumes that man is imperfect and life possibly absurd. The fiction of the former Communist Europe was based on a very different assumption, one that seems naïve and old-fashioned in its collective optimism to readers in the disillusioned democracies. It is to be noted that in the erstwhile Soviet Union aesthetic evaluation of fiction was replaced by ideological judgment. Accordingly, the works of the popular British writer A.J. Cronin, since they seem to depict personal tragedy as an emanation of capitalistic infamy, were rated higher than those of Conrad, James, and their peers.

Entertainment or escape. In a period that takes for granted that the written word should be "committed"—to the exposure of social wrong or the propagation of progressive ideologies—novelists who seek merely to take the reader out of his dull or oppressive daily life are not highly regarded, except by that reading public that has never expected a book to be anything more than a diversion. Nevertheless, the provision of laughter and dreams has been for many centuries a legitimate literary occupation. It can be condemned by serious devotees of literature only if it falsifies life through oversimplification and tends to corrupt its readers into belief that reality is as the author presents it. The novelettes once beloved of mill girls and domestic servants, in which the beggar maid was elevated to queenhood by a king of high finance, were a mere narcotic, a sort of enervating opium of the oppressed; the encouragement of such subliterate might well be one of the devices of social oppression. Adventure stories and spy novels may have a healthy enough astringency, and the very preposterousness of some adventures can be a safeguard against any impressionable young reader's neglecting the claims of real life to dream of becoming a secret agent. The subject matter of some humorous novels—such as the effete British aristocracy created by P.G. Wodehouse, which is no longer in existence if it ever was—can never be identified with a real human society; the dream is accepted as a dream. The same may be said of Evelyn Waugh's early novels—such as *Decline and Fall* (1928) and *Vile Bodies* (1930)—but these are raised

Realistic
pessimism

The
symbol-
ism of
names

above mere entertainment by touching, almost incidentally, on real human issues (the relation of the innocent to a circumambient malevolence is a persistent theme in all Waugh's writing).

Any reader of fiction has a right to an occasional escape from the dullness or misery of his existence, but he has the critical duty of finding the best modes of escape—in the most efficiently engineered detective or adventure stories, in humour that is more than sentimental buffoonery, in dreams of love that are not mere pornography. The fiction of entertainment and escape frequently sets itself higher literary standards than novels with a profound social or philosophical purpose. Books like John Buchan's *Thirty-nine Steps* (1915), Graham Greene's *Travels with My Aunt* (1969), Dashiell Hammett's *Maltese Falcon* (1930), and Raymond Chandler's *Big Sleep* (1939) are distinguished pieces of writing that, while diverting and enthralling, keep a hold on the realities of human character. Ultimately, all good fiction is entertainment, and, if it instructs or enlightens, it does so best through enchanting the reader.

Propaganda. The desire to make the reader initiate certain acts—social, religious, or political—is the essence of all propaganda, and, though it does not always accord well with art, the propagandist purpose has often found its way into novels whose prime value is an aesthetic one. The *Nicholas Nickleby* (1839) of Charles Dickens attacked the abuses of schools to some purpose, as his *Oliver Twist* (1838) drew attention to the horrors of poorhouses and his *Bleak House* (1853) to the abuses of the law of chancery. The weakness of propaganda in fiction is that it loses its value when the wrongs it exposes are righted, so that the more successful a propagandist novel is, the briefer the life it can be expected to enjoy. The genius of Dickens lay in his ability to transcend merely topical issues through the vitality with which he presented them, so that his contemporary disclosures take on a timeless human validity—chiefly through the power of their drama, character, and rhetoric.

The pure propagandist novel—which Dickens was incapable of writing—quickly becomes dated. The “social” novels of H.G. Wells, which propounded a rational mode of life and even blueprinted utopias, were very quickly exploded by the conviction of man's irredeemable irrationality that World War I initiated and World War II corroborated, a conviction the author himself came to share toward the end of his life. But the early scientific romances of Wells remain vital and are seen to have been prophetic. Most of the fiction of the former Soviet Union, which either glorified the regime or refrained from criticizing it, was dull and unreal, and the same can be said of Communist fiction elsewhere. Propaganda too frequently ignores man as a totality, concentrating on him aspectively—in terms of politics or sectarian religion. When a didactic attack on a system, as in Harriet Beecher Stowe's attack on slavery in the United States in *Uncle Tom's Cabin* (1852), seems to go beyond mere propaganda, it is because the writer makes the reader aware of wrongs and injustices that are woven into the permanent human fabric. The reader's response may be a modification of his own sensibility, not an immediate desire for action, and this is one of the legitimate effects of serious fiction. The propagandist Dickens calls for the immediate righting of wrongs, but the novelist Dickens says, mainly through implication, that all men—not just schoolmasters and state hirelings—should become more humane. If it is possible to speak of art as possessing a teaching purpose, this is perhaps its only lesson.

Reportage. The division in the novelist's mind is between his view of his art as a contrivance, like a Fabergé watch, and his view of it as a record of real life. The versatile English writer Daniel Defoe, on the evidence of such novels as his *Journal of the Plague Year* (1722), a recreation of the London plague of 1665, believed that art or contrivance had the lesser claim and proceeded to present his account of events of which he had had no direct experience in the form of plain journalistic reportage. This book, like his *Robinson Crusoe* (1719) and *Moll Flanders* (1722), is more contrived and cunning than it appears, and the hurried, unshaped narrative is the product of

careful preparation and selective ordering. His example, which could have been a very fruitful one, was not much followed until the 20th century, when the events of the real world became more terrifying and marvellous than anything the novelist could invent and seemed to ask for that full imaginative treatment that only the novelist's craft can give.

In contemporary American literature, John Hersey's *Hiroshima* (1946), though it recorded the actual results of the nuclear attack on the Japanese city in 1945, did so in terms of human immediacies, not scientific or demagogic abstractions, and this approach is essentially novelistic. Truman Capote's *In Cold Blood* (1966) took the facts of a multiple murder in the Midwest of the United States and presented them with the force, reality, tone, and (occasionally) overintense writing that distinguish his genuine fiction. Norman Mailer, in *The Armies of the Night* (1968), recorded, in great personal detail but in a third-person narration, his part in a citizens' protest march on Washington, D.C. It would seem that Mailer's talent lies in his ability to merge the art of fiction and the craft of reportage, and his *Of a Fire on the Moon* (1970), which deals with the American lunar project, reads like an episode in an emergent *roman-fleuve* of which Mailer is the central character.

The presentation of factual material as art is the purpose of such thinly disguised biographies as Somerset Maugham's *Moon and Sixpence* (1919), undisguised biographies fleshed out with supposition and imagination like Helen Waddell's *Peter Abelard* (1933), and many autobiographies served up—out of fear of libel or of dullness—as novels. Conversely, invented material may take on the lineaments of journalistic actuality through the employment of a Defoe technique of flat understatement. This is the way of such science fiction as Michael Crichton's *Andromeda Strain* (1969), which uses sketch maps, computer projections, and simulated typewritten reports.

Agent of change in language and thought. Novelists, being neither poets nor philosophers, rarely originate modes of thinking and expression. Poets such as Chaucer and Shakespeare have had much to do with the making of the English language, and Byron was responsible for the articulation of the new romantic sensibility in it in the early 19th century. Books like the Bible, Karl Marx's *Das Kapital*, and Adolf Hitler's *Mein Kampf* may underlie permanent or transient cultures, but it is hard to find, except in the early Romantic period, a novelist capable of arousing new attitudes to life (as opposed to aspects of the social order) and forging the vocabulary of such attitudes.

With the 18th-century precursors of Romanticism—Sentiment notably Richardson, Sterne, and Rousseau—the notion of sentiment entered the European consciousness. Rousseau's *Nouvelle Héloïse* fired a new attitude toward love—more highly emotional than ever before—as his *Émile* (1762) changed educated views on how to bring up children. The romantic wave in Germany, with Goethe's *Sorrows of Young Werther* (1774) and the works of Jean-Paul Richter a generation later, similarly aroused modes of feeling that rejected the rational constraints of the 18th century. Nor can the influence of Sir Walter Scott's novels be neglected, both on Europe and on the American South (where Mark Twain thought it had had a deplorable effect). With Scott came new forms of regional sentiment, based on a romantic reading of history.

It is rarely, however, that a novelist makes a profound mark on a national language, as opposed to a regional dialect (to which, by using it for a literary end, he may impart a fresh dignity). It is conceivable that Alessandro Manzoni's *I promessi sposi* (1825–27; *The Betrothed*), often called the greatest modern Italian novel, gave 19th-century Italian intellectuals some notion of a viable modern prose style in an Italian that might be termed “national,” but even this is a large claim. Günter Grass, in post-Hitler Germany, sought to revivify a language that had been corrupted by the Nazis; he threw whole dictionaries at his readers in the hope that new freedom, fantasy, and exactness in the use of words might influence the publicists, politicians, and teachers in the direction of a new liberalism of thought and expression.

The ephemeral quality of propaganda

It is difficult to say whether the French Existentialists, such as Sartre and Albert Camus, have influenced their age primarily through their fiction or their philosophical writings. Certainly, Sartre's early novel *Nausea* (1938) established unforgettable images of the key terms of his philosophy, which has haunted a whole generation, as Camus's novel *The Stranger* (1942) created for all time the lineaments of "Existential man." In the same way, the English writer George Orwell's *Nineteen Eighty-four* (1949) incarnated brilliantly the nature of the political choices that are open to 20th-century humanity, and, with terms like "Big Brother" (i.e., the leader of an authoritarian state) and "doublethink" (belief in contradictory ideas simultaneously), modified the political vocabulary. But no novelist's influence can compare to that of the poet's, who can give a language a soul and define, as Shakespeare and Dante did, the scope of a culture.

Expression of the spirit of its age. The novelist, like the poet, can make the inchoate thoughts and feelings of a society come to articulation through the exact and imaginative use of language and symbol. In this sense, his work seems to precede the diffusion of new ideas and attitudes and to be the agent of change. But it is hard to draw a line between this function and that of expressing an existing climate of sensibility. Usually the nature of a historical period—that spirit known in German as the *Zeitgeist*—can be understood only in long retrospect, and it is then that the novelist can provide its best summation. The sickness of the Germany that produced Hitler had to wait some time for fictional diagnosis in such works as Thomas Mann's *Doctor Faustus* (1947) and, later, Günter Grass's *Tin Drum* (1959). Evelyn Waugh waited several years before beginning, in the trilogy *Sword of Honour*, to depict that moral decline of English society that started to manifest itself in World War II, the conduct of which was both a cause and a symptom of the decay of traditional notions of honour and justice.

The novel can certainly be used as a tool for the better understanding of a departed age. The period following World War I had been caught forever in Hemingway's *Sun Also Rises* (1926; called *Fiesta* in England), F. Scott Fitzgerald's novels and short stories about the so-called Jazz Age, the *Antic Hay* (1923) and *Point Counter Point* (1928) of Aldous Huxley, and D.H. Lawrence's *Aaron's Rod* (1922) and *Kangaroo* (1923). The spirit of the English 18th century, during which social, political, and religious ideas associated with rising middle classes conflicted with the old Anglican Tory rigidities, is better understood through reading Smollett and Fielding than by taking the cerebral elegance of Pope and his followers as the typical expression of the period.

Similarly, the unrest and bewilderment of the young in the period after World War II still speak in novels like J.D. Salinger's *Catcher in the Rye* (1951) and Kingsley Amis' *Lucky Jim* (1954). It is notable that with novels like these—and the beat-generation books of Jack Kerouac; the American-Jewish novels of Saul Bellow, Bernard Malamud, and Philip Roth; and the black novels of Ralph Ellison and James Baldwin—it is a segmented spirit that is expressed, the spirit of an age group, social group, or racial group, and not the spirit of an entire society in a particular phase of history. But probably a *Zeitgeist* has always been the emanation of a minority, the majority being generally silent. The 20th century seems, from this point of view, to be richer in vocal minorities than any other period in history.

Creator of life-style and arbiter of taste. Novels have been known to influence, though perhaps not very greatly, modes of social behaviour and even, among the very impressionable, conceptions of personal identity. But more young men have seen themselves as Hamlet or Childe Harold than as Julien Sorel, the protagonist of Stendhal's novel *The Red and the Black* (1830), or the sorrowing Werther. Richardson's novel may popularize Pamela, or Galsworthy's *Forsyte Saga* (1906–22) Jon, as a baptismal name, but it rarely makes a deeper impression on the mode of life of literate families. On the other hand, the capacity of Oscar Wilde's *Picture of Dorian Gray* (1891) to influence young men in the direction of sybaritic amoral-

ity, or of D.H. Lawrence's *Lady Chatterley's Lover* (1928) to engender a freer attitude to sex, has never been assessed adequately. With the lower middle class reading public, the effect of devouring *The Forsyte Saga* was to engender gentelisms—cucumber sandwiches for tea, supper renamed dinner—rather than to learn that book's sombre lesson about the decline of the old class structure. Similarly, the ladies who read Scott in the early 19th century were led to barbarous ornaments and tastefully arranged folk songs.

Fiction has to be translated into one of the dramatic media—stage, film, or television—before it can begin to exert a large influence. *Tom Jones* as a film in 1963 modified table manners and coiffures and gave American visitors to Great Britain a new (and probably false) set of expectations. The stoic heroes of Hemingway, given to drink, fights, boats, and monosyllables, became influential only when they were transferred to the screen. They engendered other, lesser heroes—incorruptible private detectives, partisans brave under interrogation—who in their turn have influenced the impressionable young when seeking an identity. Ian Fleming's James Bond led to a small revolution in martini ordering. But all these influences are a matter of minor poses, and such poses are most readily available in fiction easily adapted to the mass media—which means lesser fiction. Proust, though he recorded French patrician society with painful fidelity, had little influence on it, and it is hard to think of Henry James disturbing the universe even fractionally. Films and television programs dictate taste and behaviour more than the novel ever could.

STYLE

Romanticism. The Romantic movement in European literature is usually associated with those social and philosophical trends that prepared the way for the French Revolution, which began in 1789. The somewhat subjective, anti-rational, emotional currents of romanticism transformed intellectual life in the revolutionary and Napoleonic periods and remained potent for a great part of the 19th century. In the novel, the romantic approach to life was prepared in the "sentimental" works of Richardson and Sterne and attained its first major fulfillment in the novels of Rousseau. Sir Walter Scott, in his historical novels, turned the past into a great stage for the enactment of events motivated by idealism, chivalry, and strong emotional impulse, using an artificially archaic language full of remote and magical charm. The exceptional soul—poet, patriot, idealist, madman—took the place of dully reasonable fictional heroes, such as Tom Jones, and sumptuous and mysterious settings ousted the plain town and countryside of 18th-century novels.

The romantic novel must be seen primarily as a historical phenomenon, but the romantic style and spirit, once they had been brought into being, remained powerful and attractive enough to sustain a whole subspecies of fiction. The cheapest love story can be traced back to the example of Charlotte Brontë's *Jane Eyre* (1847), or even Rousseau's earlier *Nouvelle Héloïse*. Similarly, best-selling historical novels, even those devoid of literary merit, can find their progenitor in Scott, and science fiction in Mary Shelley's *Frankenstein* (1818), a romantic novel subtitled *The Modern Prometheus*, as well as in Jules Verne and H.G. Wells. The aim of romantic fiction is less to present a true picture of life than to arouse the emotions through a depiction of strong passions, or to fire the imagination with exotic, terrifying, or wonderful scenes and events. When it is condemned by critics, it is because it seems to falsify both life and language: the pseudopoetical enters the dialogue and *vérité* alike, and humanity is seen in only one of its aspects—that of feeling untempered with reason.

If such early romantic works as those of Scott and of the Goethe of *The Sorrows of Werther* have long lost their original impact, the romantic spirit still registers power and truth in the works of the Brontës—particularly in Emily Brontë's *Wuthering Heights*, in which the poetry is genuine and the strange instinctual world totally convincing. Twentieth-century romantic fiction records few masterpieces. Writers like Daphne du Maurier, the author of *Jamaica Inn* (1936), *Rebecca* (1938), and many others,

The lag between events and their use in novels

Fiction in other media

Decline of the romantic novel

are dismissed as mere purveyors of easy dreams. It is no more possible in the 20th century to revive the original romantic élan in literature than it is to compose music in the style of Beethoven. Despite the attempts of Lawrence Durrell to achieve a kind of decadent romantic spirit in his *Alexandria Quartet*, the strong erotic feeling, the exotic setting, the atmosphere of poetic hallucination, the pain, perversion, and elemental force seem to be contrivances, however well they fulfill the original romantic prescription.

Realism. Certain major novelists of the 19th century, particularly in France, reacted against romanticism by eliminating from their work those "softer" qualities—tenderness, idealism, chivalric passion, and the like—which seemed to them to hide the stark realities of life in a dreamlike haze. In Gustave Flaubert's works there are such romantic properties—his novel *Salammbô* (1862), for instance, is a sumptuous representation of a remote pagan past—but they are there only to be punctured with realistic irony. On one level, his *Madame Bovary* may be taken as a kind of parable of the punishment that fate metes out to the romantic dreamer, and it is the more telling because Flaubert recognized a strong romantic vein in himself: "Madame Bovary, c'est moi" ("Madame Bovary is myself"). Stendhal and Balzac, on the other hand, admit no dreams and present life in a grim nakedness without poetic drapery.

Balzac's mammoth fictional work—the 20-year succession of novels and stories he published under the collective title *La Comédie humaine* (*The Human Comedy*)—and Stendhal's novels of the same period, *The Red and the Black* (1830) and *The Charterhouse of Parma* (1839), spare the reader nothing of those baser instincts in man and society that militate against, and eventually conquer, many human aspirations. Rejecting romanticism so energetically, however, they swing to an extreme that makes "realism" a synonym for unrelenting pessimism. Little comes right for the just or the weak, and base human nature is unqualified by even a modicum of good. But there is a kind of affirmative richness and energy about both writers that seems to belie their pessimistic thesis.

In England, George Eliot in her novel *Middlemarch* (1871–72) viewed human life grimly, with close attention to the squalor and penury of rural life. If "nature" in works by romantic poets like Wordsworth connoted a kind of divine benevolence, only the "red in tooth and claw" aspect was permitted to be seen in the novels of the realists. George Eliot does not accept any notion of Divine Providence, whether Christian or pantheistic, but her work is instinct with a powerful moral concern: her characters never sink into a deterministic morass of hopelessness, since they have free will, or the illusion of it. With Thomas Hardy, who may be termed the last of the great 19th-century novelists, the determinism is all-pervasive, and his final novel, *Jude the Obscure* (1896), represents the limit of pessimism. Behind him one is aware of the new science, initiated by the biologists Charles Darwin and T.H. Huxley, which displaces man as a free being, capable of choice, by a view of him as the product of blind mechanistic forces over which he has little control.

Realism in this sense has been a continuing impulse in the 20th-century novel, but few writers would go so far as Hardy in positing man's near-total impotence in a hostile universe, with the gods killing human creatures for their sport. Realism in the Existentialist fiction of 20th-century France, for instance, makes man not merely wretched but absurd, yet it does not diminish his power of self-realization through choice and action. Realism has frequently been put in the service of a reforming design, which implies a qualified optimism. War novels, novels about the sufferings of the oppressed (in prison, ghetto, totalitarian state), studies of human degradation that are bitter cries against man-made systems—in all of these the realistic approach is unavoidable, and realistic detail goes much further than anything in the first realists. But there is a difference in the quality of the anger the reader feels when reading the end of Hardy's *Tess of the D'Urbervilles* (1891) and that generated by Upton Sinclair's *Jungle* (1906) or Erich Maria Remarque's *All Quiet on the Western Front* (1929). In Hardy's novel, pessimistic de-

terminism, reducing human character to pain, frustration, and impotent anger, was—paradoxically—appropriate to an age that knew no major cataclysms or oppressions. The novels of Sinclair and Remarque reflect the 20th century, which saw the origin of all wrong in the human will, and set on a program of diagnosis and reform.

Naturalism. The naturalistic novel is a development out of realism, and it is, again, in France that its first practitioners are to be found, with Émile Zola leading. It is difficult to separate the two categories, but naturalism seems characterized not only by a pessimistic determinism but also by a more thoroughgoing attention to the physical and biological aspects of human existence. Man is less a soul aspiring upward to its divine source than a product of natural forces, as well as genetic and social influences, and the novelist's task is to present the physical essence of man and his environment. The taste of Balzac's and Stendhal's audiences was not easily able to accommodate itself to utter frankness about the basic processes of life, and the naturalists had to struggle against prejudice, and often censorship, before their literary candour was able to prevail. The 20th century takes the naturalistic approach for granted, but it is more concerned with a technique of presentation than with the somewhat mechanistic philosophy of Zola and his followers.

Naturalism received an impetus after World War I, when novelists felt they had a duty to depict the filth, suffering, and degradation of the soldier's life, without euphemism or circumlocution. Joyce's *Ulysses*, when it appeared in 1922, was the first novel to seek to justify total physical candour in terms of its artistic, as opposed to moral, aim—which was to depict with almost scientific objectivity every aspect of an ordinary urban day. Though Joyce had read Zola, he seems to invoke the spirit of a very much earlier naturalistic writer—the ribald French author of the 16th century, François Rabelais—and this is in keeping with the Catholic tradition that Joyce represents. Zola, of course, was an atheist.

It would have been a sin against his aesthetic canons for Joyce to have shown Leopold Bloom—the protagonist of *Ulysses*—eating breakfast or taking a bath and yet not defecating or masturbating. The technique of the interior monologue, which presented the unedited flow of a character's unspoken thought and emotion, also called for the utmost frankness in dealing with natural functions and urges. Joyce, it is now recognized, had no prurient or scatological intention; his concern was with showing life as it is (without any of the didactic purpose of Zola), and this entailed the presentation of lust, perversion, and blasphemy as much as any of the traditionally acceptable human functions.

The naturalistic novelists have had their social and legal problems—obscenity indictments, confiscation, emasculation by timid publishers—but the cause was ultimately won, at least in Great Britain and the United States, where there are few limits placed on the contemporary novelist's proclaimed right to be true to nature. In comparison with much contemporary fiction the pioneer work of Zola seems positively reticent.

Impressionism. The desire to present life with frank objectivity led certain early 20th-century novelists to question the validity of long-accepted narrative conventions. If truth was the novelist's aim, then the tradition of the omniscient narrator would have to go, to be replaced by one in which a fallible, partially ignorant character—one involved in the story and hence himself subject to the objective or naturalistic approach—recounted what he saw and heard. But the Impressionist painters of late 19th-century France had proclaimed a revision of the whole seeing process: they distinguished between what the observer assumed he was observing and what he actually observed. That cerebral editing which turned visual data into objects of geometric solidity had no place in Impressionist painting; the visible world became less definite, more fluid, resolving into light and colour.

The German novelists Thomas Mann and Hermann Hesse, moving from the realist tradition, which concentrated on closely notated detail in the exterior world, sought the lightness and clarity of a more elliptical style, and were

The pessimistic view

Realism as reform

The influence of the painters

proclaimed Impressionists. But in England Ford Madox Ford went much further in breaking down the imagined rigidities of the space-time continuum, liquidating step-by-step temporal progression and making the visual world shimmer, dissolve, reconstitute itself. In Ford's tetralogy *Parade's End* (1924–28), the reader moves freely within the time continuum, as if it were spatial, and the total picture is perceived through an accumulation of fragmentary impressions. Ford's masterpiece, *The Good Soldier*, pushes the technique to its limit: the narrator tells his story with no special dispensation to see or understand more than a fallible being can, and, in his reminiscences, he fragments whole sequences of events as he ranges freely through time (such freedom had traditionally been regarded as a weakness, a symptom of the disease of inattention).

In the approach to dialogue manifested in a book that Ford wrote jointly with Conrad—*The Inheritors* (1901)—a particular aspect of literary impressionism may be seen whose suggestiveness has been ignored by other modern novelists. As the brain imposes its own logical patterns on the phenomena of the visual world, so it is given to editing into clarity and conciseness the halting utterances of real-life speech; the characters of most novels are impossibly articulate. Ford and Conrad attempted to present speech as it is actually spoken, with many of the meaningful solidities implied rather than stated. The result is sometimes exasperating, but only as real-life conversation frequently is.

The interior monologue, which similarly resists editing, may be regarded as a development of this technique. To show pre-articulatory thought, feeling, and sensuous perception unordered into a rational or "literary" sequence is an impressionistic device that, beginning in Édouard Dujardin's minor novel *Les Lauriers sont coupés* (1888; *We'll to the Woods No More*), served fiction of high importance, from Dorothy Richardson, Joyce, and Virginia Woolf to William Faulkner and Samuel Beckett.

Novelists like Ronald Firbank and Evelyn Waugh (who studied painting and was a competent draftsman) learned, in a more general sense, how to follow the examples of the Impressionist and Postimpressionist painters in their fiction. A spare brilliance of observation, like those paintings in which a whole scene is suggested through carefully selected points of colour, replaced that careful delineation of a whole face, or inventorying of a whole room, that had been the way of Balzac and other realists. In four or five brief lines of dialogue Waugh can convey as much as the 19th-century novelists did in as many pages.

Expressionism. Expressionism was a German movement that found its most congenial media in painting and drama. The artist's aim was to express, or convey the essence of, a particular theme, to the exclusion of such secondary considerations as fidelity to real life. The typical Expressionist play, by Bertolt Brecht, for example, concerns itself with a social or political idea that is hurled at the audience through every possible stage device—symbols, music, cinematic insertions, choral speech, dance. Human character is less important than the idea of humanity, and probability of action in the old realist sense is the least of the dramatist's concerns. The emotional atmosphere is high-pitched, even ecstatic, and the tone is more appropriate to propaganda than to art. Expressionistic technique, as the plays of Brecht prove, was an admirable means of conveying a Communist program, and it was in the service of such a program that John Dos Passos, in the trilogy of novels *U.S.A.* (1937), used literary devices analogous to the dramatic ones of Brecht—headlines, tabloid biographies, popular songs, lyric soliloquies, and the like.

But the Austro-Czech Franz Kafka, the greatest of the Expressionist novelists, sought to convey what may crudely be termed man's alienation from his world in terms that admit of no political interpretation. Joseph K., the hero of Kafka's novel *The Trial* (1925), is accused of a nameless crime, he seeks to arm himself with the apparatus of a defense, and he is finally executed—stabbed with the utmost courtesy by two men in a lonely place. The hallucinatory atmosphere of that novel, as also of his novel *The Castle* (1926), is appropriate to nightmare, and indeed Kafka's work has been taken by many as an imaginative forecast

of the nightmare through which Europe was compelled to live during the Hitler regime. But its significance is more subtle and universal: one of the elements is original sin and another filial guilt. In the story *The Metamorphosis* (1915) a young man changes into an enormous insect, and the nightmare of alienation can go no further.

Kafka's influence has been considerable. Perhaps his most distinguished follower is the English writer Rex Warner, whose *Wild Goose Chase* (1937) and *Aerodrome* (1941) use fantasy, symbol, and improbable action for an end that is both Marxist and Freudian: the filial guilt, however, seems to be taken directly from Kafka, with an innocent hero caught in a monstrously oppressive web that is both the totalitarian state and paternal tyranny. More recently, the American writer William Burroughs has developed his own Expressionistic techniques in *The Naked Lunch* (1959), which is concerned with the alienation from society of the drug addict. His later novels *Nova Express* (1964) and *The Ticket That Exploded* (1962) use obscene fantasy to present a kind of metaphysical struggle between free spirit and enslaved flesh, evidently an extrapolation of the earlier drug theme. Burroughs is a didactic novelist, and didacticism functions best in a fictional ambience that rejects the complexities of character and real-life action.

Avant-gardism. Many innovations in fiction can be classified under headings already considered. Even so revolutionary a work as Joyce's *Finnegans Wake* represents an attempt to show the true nature of a dream; this can be regarded as a kind of Impressionism pushed so far that it looks like Surrealism. The brief novels of Samuel Beckett (which, as they aim to demonstrate the inadequacy of language to express the human condition, become progressively more brief) seem to have a kind of Expressionist derivation, since everything in them is subordinated to a central image of man as a totally deprived creature, resentful of a God he does not believe in. The French anti-novel, dethroning man as a primary concern of fiction, perhaps represents the only true break with traditional technique that the 20th-century novel has seen.

Dissatisfaction not only with the content of the traditional novel but with the manner in which readers have been schooled to approach it has led the contemporary French novelist Michel Butor, in *Mobile*, to present his material in the form of a small encyclopaedia, so that the reader finds his directions obliquely, through an alphabetic taxonomy and not through the logic of sequential events. Nabokov, in *Pale Fire* (1962), gives the reader a poem of 999 lines and critical apparatus assembled by a madman; again the old sense of direction (beginning at the beginning and going on to the end) has been liquidated, yet *Pale Fire* is a true and highly intelligible novel. In England, B.S. Johnson published similar "false-directional" novels, though the influence of Sterne makes them seem accessible, even cozily traditional. One of Johnson's books is marketed as a bundle of disjoint chapters—which may thus be dealt aleatorially and read in any order.

Available avant-garde techniques are innumerable, though not all of them are salable. There is the device of counterpointing a main narrative with a story in footnotes, which eventually rises like water and floods the other. A novel has been written, though not published, in which the words are set (rather like the mouse's tail or tale in *Alice in Wonderland*) to represent graphically the physical objects in the narrative. Burroughs has experimented with a tricolunar technique, in which three parallel narratives demand the reader's attention. But the writers like Borges and Nabokov go beyond mere technical innovation: they ask for a reconsideration of the very essence of fiction. In one of his *ficciones*, Borges strips from the reader even the final illusion that he is reading a story, for the story is made to dissolve, the artist evidently losing faith in his own artifact. Novels, as both Borges and Nabokov show, can turn into poems or philosophical essays, but they cannot, while remaining literature, turn into compositions disclaiming all interest in the world of feeling, thought, and sense. The novelist can do anything he pleases with his art so long as he interprets, or even just presents, a world that the reader recognizes as existing, or capable of existing, or capable of being dreamed of as existing.

Expression-
ism as
alien-
ation

Technical
innovation

TYPES OF NOVEL

Historical. For the hack novelist, to whom speedy output is more important than art, thought, and originality, history provides ready-made plots and characters. A novel on Alexander the Great or Joan of Arc can be as flimsy and superficial as any schoolgirl romance. But historical themes, to which may be added prehistoric or mythical ones, have inspired the greatest novelists, as Tolstoy's *War and Peace* and Stendhal's *Charterhouse of Parma* reveal. In the 20th century, distinguished historical novels such as Arthur Koestler's *The Gladiators* (1939), Robert Graves's *I, Claudius* (1934), Zoë Oldenbourg's *Destiny of Fire* (1960), and Mary Renault's *The King Must Die* (1958) exemplify an important function of the fictional imagination—to interpret remote events in human and particular terms, to transform documentary fact, with the assistance of imaginative conjecture, into immediate sensuous and emotional experience.

There is a kind of historical novel, little more than a charade, which frequently has a popular appeal because of a common belief that the past is richer, bloodier, and more erotic than the present. Such novels, which include such immensely popular works as those of Georgette Heyer, or Baroness Orczy's Scarlet Pimpernel stories in England in the early 20th century, and *Forever Amber* (1944) by Kathleen Winsor in the United States, may use the trappings of history but, because there is no real assimilation of the past into the imagination, the result must be a mere costume ball. On the other hand, the American novelist John Barth showed in *The Sot-Weed Factor* (1960) that mock historical scholarship—preposterous events served up with parodic pomposity—could constitute a viable, and not necessarily farcical, approach to the past. Barth's history is cheerfully suspect, but his sense of historical perspective is genuine.

It is in the technical conservatism of most European historical novels that the serious student of fiction finds cause to relegate the category to a secondary place. Few practitioners of the form seem prepared to learn from any writer later than Scott, though Virginia Woolf—in *Orlando* (1928) and *Between the Acts* (1941)—made bold attempts to squeeze vast tracts of historical time into a small space and thus make them as fictionally manageable as the events of a single day. And John Dos Passos' *U.S.A.*, which can be taken as a historical study of a phase in America's development, is a reminder that experiment is not incompatible with the sweep and amplitude that great historical themes can bring to the novel.

Picaresque. In Spain, the novel about the rogue or *pícaro* was a recognized form, and such English novels as Defoe's *The Fortunate Mistress* (1724) can be regarded as picaresque in the etymological sense. But the term has come to connote as much the episodic nature of the original species as the dynamic of roguery. Fielding's *Tom Jones*, whose hero is a bastard, amoral, and very nearly gallows-meat, has been called picaresque, and the *Pickwick Papers* of Dickens—whose eponym is a respectable and even childishly ingenuous scholar—can be accommodated in the category.

The requirements for a picaresque novel are apparently length, loosely linked episodes almost complete in themselves, intrigue, fights, amorous adventure, and such optional items as stories within the main narrative, songs, poems, or moral homilies. Perhaps inevitably, with such a structure or lack of it, the driving force must come from a wild or roguish rejection of the settled bourgeois life, a desire for the open road, with adventures in inn bedrooms and meetings with questionable wanderers. In the modern period, Saul Bellow's *Adventures of Augie March* (1953) and Jack Kerouac's *Dharma Bums* (1959) have something of the right episodic, wandering, free, questing character. But in an age that lacks the unquestioning acceptance of traditional morality against which the old picaresque heroes played out their villainous lives, it is not easy to revive the *novela picaresca* as the anonymous author of *Lazarillo de Tormes* (1554) conceived it, or as such lesser Spanish writers of the beginning of the 17th century as Mateo Alemán, Vicente Espinel, and Luis Vélez de Guevara developed it. The modern criminal wars with the

police rather than with society, and his career is one of closed and narrow techniques, not compatible with the gay abandon of the true *pícaro*.

Sentimental. The term sentimental, in its mid-18th-century usage, signified refined or elevated feeling, and it is in this sense that it must be understood in Laurence Sterne's *Sentimental Journey* (1768). Richardson's *Pamela* (1740) and Rousseau's *Nouvelle Héloïse* (1761) are sentimental in that they exhibit a passionate attachment between the sexes that rises above the merely physical. The vogue of the sentimental love novel was one of the features of the Romantic movement, and the form maintained a certain moving dignity despite a tendency to excessive emotional posturing. The germs of mawkishness are clearly present in Sterne's *Tristram Shandy* (1760–67), though offset by a diluted Rabelaisianism and a certain cerebral quality. The debasement by which the term sentimental came to denote a self-indulgence in superficial emotions occurred in the Victorian era, under the influence of sanctimony, religiosity, and a large commercial demand for bourgeois fiction. Sentimental novels of the 19th and 20th centuries are characterized by an invertebrate emotionalism and a deliberately lachrymal appeal. Neither Dickens nor Thackeray was immune to the temptations of sentimentality—as is instanced by their treatment of deathbed scenes. The reported death of Tiny Tim in *A Christmas Carol* (1843) is an example of Dickens' ability to provoke two tearful responses from the one situation—one of sorrow at a young death, the other of relief at the discovery that the death never occurred. Despite such patches of emotional excess, Dickens cannot really be termed a sentimental novelist. Such a designation must be reserved for writers like Mrs. Henry Wood, the author of *East Lynne* (1861). That the sentimental novel is capable of appeal even in the Atomic Age is shown by the success of *Love Story* (1970), by Erich Segal. That this is the work of a Yale professor of classics seems to indicate either that not even intellectuals disdain sentimental appeal or that tearjerking is a process to be indulged in coldly and even cynically. Stock emotions are always easily aroused through stock devices, but both the aim and the technique are generally eschewed by serious writers.

Gothic. The first Gothic fiction appeared with works like Horace Walpole's *Castle of Otranto* (1765) and Matthew Gregory Lewis' *Monk* (1796), which countered 18th-century "rationalism" with scenes of mystery, horror, and wonder. Gothic (the spelling "Gothick" better conveys the contemporary flavour) was a designation derived from architecture, and it carried—in opposition to the Italianate style of neoclassical building more appropriate to the Augustan Age—connotations of rough and primitive grandeur. The atmosphere of a Gothic novel was expected to be dark, tempestuous, ghostly, full of madness, outrage, superstition, and the spirit of revenge. Mary Shelley's *Frankenstein*, which maintains its original popularity and even notoriety, has in addition the traditional Gothic ingredients, with its weird God-defying experiments, its eldritch shrieks, and, above all, its monster. Edgar Allan Poe developed the Gothic style brilliantly in the United States, and he has been a considerable influence. A good deal of early science fiction, like H.G. Wells's *Island of Doctor Moreau* (1896), seems to spring out of the Gothic movement, and the Gothic atmosphere has been seriously cultivated in England in the later novels of Iris Murdoch and in the Gormenghast sequence beginning in 1946 of Mervyn Peake. It is noteworthy that Gothic fiction has always been approached in a spirit of deliberate suspension of the normal canons of taste. Like a circus trick, a piece of Gothic fiction asks to be considered as ingenious entertainment; the pity and terror are not aspects of a cathartic process but transient emotions to be, somewhat perversely, enjoyed for their own sake.

Psychological. The psychological novel first appeared in 17th-century France, with Madame de La Fayette's *Princesse de Clèves* (1678), and the category was consolidated by works like the Abbé Prévost's *Manon Lescaut* (1731) in the century following. More primitive fiction had been characterized by a proliferation of action and incidental characters; the psychological novel limited it-

Ingredients of the typical Gothic novel

self to a few characters whose motives for action could be examined and analyzed. In England, the psychological novel did not appear until the Victorian era, when George Eliot became its first great exponent. It has been assumed since then that the serious novelist's prime concern is the workings of the human mind, and hence much of the greatest fiction must be termed psychological. Dostoyevsky's *Crime and Punishment* deals less with the ethical significance of a murder than with the soul of the murderer; Flaubert's interest in Emma Bovary has less to do with the consequences of her mode of life in terms of nemesic logic than with the patterns of her mind; in *Anna Karenina*, Tolstoy presents a large-scale obsessive study of feminine psychology that is almost excruciating in its relentless probing. The novels of Henry James are psychological in that the crucial events occur in the souls of the protagonists, and it was perhaps James more than any serious novelist before or since who convinced frivolous novel-readers that the "psychological approach" guarantees a lack of action and excitement.

The theories of Sigmund Freud are credited as the source of the psychoanalytical novel. Freud was anticipated, however, by Shakespeare (in, for example, his treatment of Lady Macbeth's somnambulistic guilt). Two 20th-century novelists of great psychological insight—Joyce and Nabokov—professed a disdain for Freud. To write a novel with close attention to the Freudian or Jungian techniques of analysis does not necessarily produce new prodigies of psychological revelation; Oedipus and Electra complexes have become commonplaces of superficial novels and films. The great disclosures about human motivation have been achieved more by the intuition and introspection of novelists and dramatists than by the more systematic work of the clinicians.

The novel of manners. To make fiction out of the observation of social behaviour is sometimes regarded as less worthy than to produce novels that excavate the human mind. And yet the social gestures known as manners, however superficial they appear to be, are indices of a collective soul and merit the close attention of the novelist and reader alike. The works of Jane Austen concern themselves almost exclusively with the social surface of a fairly narrow world, and yet she has never been accused of a lack of profundity. A society in which behaviour is codified, language restricted to impersonal formulas, and the expression of feeling muted, is the province of the novel of manners, and such fiction may be produced as readily in the 20th century as in the era of Fanny Burney or Jane Austen. Such novels as Evelyn Waugh's *Handful of Dust* (1934) depend on the exact notation of the manners of a closed society, and personal tragedies are a mere temporary disturbance of collective order. Even Waugh's trilogy *Sword of Honour* is as much concerned with the minutiae of surface behaviour in an army, a very closed society, as with the causes for which that army fights. H.H. Munro ("Saki"), in *The Unbearable Bassington* (1912), an exquisite novel of manners, says more of the nature of Edwardian society than many a more earnest work. It is conceivable that one of the novelist's duties to posterity is to inform it of the surface quality of the society that produced him; the great psychological profundities are eternal, manners are ephemeral and have to be caught. Finally, the novel of manners may be taken as an artistic symbol of a social order that feels itself to be secure.

Epistolary. The novels of Samuel Richardson arose out of his pedagogic vocation, which arose out of his trade of printer—the compilation of manuals of letter-writing technique for young ladies. His age regarded letter writing as an art on which could be expended the literary care appropriate to the essay or to fiction, and, for Richardson, the creation of epistolary novels entailed a mere step from the actual world into that of the imagination. His *Pamela* (1740) and *Clarissa* (1748) won phenomenal success and were imitated all over Europe, and the epistolary novel—with its free outpouring of the heart—was an aspect of early romanticism. In the 19th century, when the letter-writing art had not yet fallen into desuetude, it was possible for Wilkie Collins to tell the mystery story of *The Moonstone* (1868) in the form of an exchange of letters,

but it would be hard to conceive of a detective novel using such a device in the 20th century, when the well-wrought letter is considered artificial. Attempts to revive the form have not been successful, and Christopher Isherwood's *Meeting by the River* (1967), which has a profoundly serious theme of religious conversion, seems to fail because of the excessive informality and chattiness of the letters in which the story is told. The 20th century's substitute for the long letter is the transcribed tape recording—more, as Beckett's play *Krapp's Last Tape* indicates, a device for expressing alienation than a tool of dialectic. But it shares with the Richardsonian epistle the power of seeming to grant direct communication with a fictional character, with no apparent intervention on the part of the true author.

Pastoral. Fiction that presents rural life as an idyllic condition, with exquisitely clean shepherdesses and sheep immune to foot-rot, is of very ancient descent. Longus' *Daphnis and Chloe*, written in Greek in the 2nd or 3rd century AD, was the remote progenitor of such Elizabethan pastoral romances as Sir Philip Sidney's *Arcadia* (1590) and Thomas Lodge's *Rosalynde* (1590), the source book for Shakespeare's *As You Like It*. The *Paul et Virginie* of Bernardin de St. Pierre (1787), which was immensely popular in its day, seems to spring less from the pastoral utopian convention than from the dawning Romanticism that saw in a state of nature only goodness and innocence. Still, the image of a rural Eden is a persistent one in Western culture, whatever the philosophy behind it, and there are elements of this vision even in D.H. Lawrence's *Rainbow* (1915) and, however improbable this may seem, in his *Lady Chatterley's Lover* (1928). The more realistic and ironic pictures of the pastoral life, with poverty and pig dung, beginning with George Crabbe's late-18th-century narrative poems, continuing in George Eliot, reaching sour fruition in Thomas Hardy, are usually the work of people who know the country well, while the rural idyll is properly a townsman's dream. The increasing stresses of urban life make the country vision a theme still available to serious fiction, as even a work as sophisticated as Saul Bellow's *Herzog* (1964) seems to show. But, since Stella Gibbons' satire *Cold Comfort Farm* (1932), it has been difficult for any British novelist to take seriously pastoral lyricism.

Apprenticeship. The *Bildungsroman*, or novel about upbringing and education, seems to have its beginnings in Goethe's work, *Wilhelm Meisters Lehrjahre* (1796), which is about the processes by which a sensitive soul discovers its identity and its role in the big world. A story of the emergence of a personality and a talent, with its implicit motifs of struggle, conflict, suffering, and success, has an inevitable appeal for the novelist; many first novels are autobiographical and attempt to generalize the author's own adolescent experiences into a kind of universal symbol of the growing and learning processes. Charles Dickens embodies a whole *Bildungsroman* in works like *David Copperfield* (1850) and *Great Expectations* (1861), but allows the emerged ego of the hero to be absorbed into the adult world, so that he is the character that is least remembered. H.G. Wells, influenced by Dickens but vitally concerned with education because of his commitment to socialist or utopian programs, looks at the agonies of the growing process from the viewpoint of an achieved utopia in *The Dream* (1924) and, in *Joan and Peter* (1918), concentrates on the search for the right modes of apprenticeship to the complexities of modern life.

The school story established itself in England as a form capable of popularization in children's magazines, chiefly because of the glamour of elite systems of education as first shown in Thomas Hughes's *Tom Brown's School Days* (1857), which is set at Rugby. In France, *Le Grand Meaulnes* (1913) of Alain-Fournier is the great exemplar of the school novel. The studies of struggling youth presented by Hermann Hesse became, after his death in 1962, part of an American campus cult indicating the desire of the serious young to find literary symbols for their own growing problems.

Samuel Butler's *Way of All Flesh*, which was written by 1885 but not published until 1903, remains one of the

The movement away from Freud

Realism in the pastoral novel

greatest examples of the modern *Bildungsroman*: philosophical and polemic as well as moving and comic, it presents the struggle of a growing soul to further, all unconsciously, the aims of evolution, and is a devastating indictment of Victorian paternal tyranny. But probably James Joyce's *Portrait of the Artist as a Young Man* (1916), which portrays the struggle of the nascent artistic temperament to overcome the repressions of family, state, and church, is the unsurpassable model of the form in the 20th century. That the learning novel may go beyond what is narrowly regarded as education is shown in two remarkable works of the 1950s—William Golding's *Lord of the Flies* (1955), which deals with the discovery of evil by a group of shipwrecked middle class boys brought up in the liberal tradition, and J.D. Salinger's *Catcher in the Rye* (1951), which concerns the attempts of an adolescent American to come to terms with the adult world in a series of brief encounters, ending with his failure and his ensuing mental illness.

Roman à clef. Real, as opposed to imaginary, human life provides so much ready-made material for the novelist that it is not surprising to find in many novels a mere thinly disguised and minimally reorganized representation of actuality. When, for the fullest appreciation of a work of fiction, it is necessary for the reader to consult the real-life personages and events that inspired it, then the work is a *roman à clef*, or novel that needs a key. In a general sense, every work of literary art requires a key or clue to the artist's preoccupations (the jail in Dickens; the mysterious tyrants in Kafka, both leading back to the author's own father), but the true *roman à clef* is more particular in its disguised references. Chaucer's "Nun's Priest's Tale" has puzzling naturalistic details that can be cleared up only by referring the poem to an assassination plot in which the Earl of Bolingbroke was involved. Swift's *Tale of a Tub* (1704), Dryden's *Absalom and Achitophel* (1681), and Orwell's *Animal Farm* (1945) make total sense only when their hidden historical content is disclosed. These, of course, are not true novels, but they serve to indicate a literary purpose that is not primarily aesthetic. Lawrence's *Aaron's Rod* requires a knowledge of the author's personal enmities, and to understand Aldous Huxley's *Point Counter Point* fully one must know, for instance, that the character of Mark Rampion is D.H. Lawrence himself and that of Denis Burlap is the critic John Middleton Murry. Proust's *À la recherche du temps perdu* becomes a richer literary experience when the author's social milieu is explored, and Joyce's *Finnegans Wake* has so many personal references that it may be called the most massive *roman à clef* ever written. The more important the *clef* becomes to full understanding, the closer the work has come to a special kind of didacticism. When it is dangerous to expose the truth directly, then the novel or narrative poem may present it obliquely. But the ultimate vitality of the work will depend on those elements in it that require no key.

Anti-novel. The movement away from the traditional novel form in France in the form of the *nouveau roman* tends to an ideal that may be called the anti-novel—a work of the fictional imagination that ignores such properties as plot, dialogue, human interest. It is impossible, however, for a human creator to create a work of art that is completely inhuman. Contemporary French writers like Alain Robbe-Grillet in *Jealousy* (1957), Nathalie Sarraute in *Tropisms* (1939) and *The Planetarium* (1959), and Michel Butor in *Passing Time* (1957) and *Degrees* (1960) wish mainly to remove the pathetic fallacy from fiction, in which the universe, which is indifferent to man, is made to throw back radar reflections of man's own emotions. Individual character is not important, and consciousness dissolves into sheer "perception." Even time is reversible, since perceptions have nothing to do with chronology, and, as Butor's *Passing Time* shows, memories can be lived backward in this sort of novel. Ultimately, the very appearance of the novel—traditionally a model of the temporal treadmill—must change; it will not be obligatory to start at page 1 and work through to the end; a novel can be entered at any point, like an encyclopaedia.

The two terms most heard in connection with the French anti-novel are *chosisme* and *tropisme*. The first, with which

Robbe-Grillet is chiefly associated, relates to the novelist's concern with things in themselves, not things as human symbols or metaphors. The second, which provided a title for Nathalie Sarraute's early novel, denotes the response of the human mind to external stimuli—a response that is general and unmodified by the apparatus of "character." It is things, the furniture of the universe, that are particular and variable; the multiplicity of human observers melts into an undifferentiable mode of response. Needless to say, there is nothing new in this epistemology as applied to the novel. It is present in Laurence Sterne (in whom French novelists have always been interested), as also in Virginia Woolf.

Such British practitioners of the anti-novel as Christine Brooke-Rose and Rayner Heppenstall (both French scholars, incidentally) are more empirical than their French counterparts. They object mainly to the falsification of the external world that was imposed on the traditional novel by the exigencies of plot and character, and they insist on notating the minutiae of the surface of life, concentrating in an unhurried fashion on every detail of its texture. A work like Heppenstall's *Connecting Door* (1962), in which the narrator-hero does not even possess a name, is totally unconcerned with action but very interested in buildings, streets, and the sound of music. This is properly a fresh approach to the materials of the traditional novel rather than a total liberation from it. Such innovations as are found in the *nouveau roman* can best show their value in their influence on traditional novelists, who may be persuaded to observe more closely and be wary of the seductions of swift action, contrived relationships, and neat resolutions.

Cult, or coterie, novels. The novel, unlike the poem, is a commercial commodity, and it lends itself less than the materials of literary magazines to that specialized appeal called coterie, intellectual or elitist. It sometimes happens that books directed at highly cultivated audiences—like *Ulysses*, *Finnegans Wake*, and Djuna Barnes's *Nightwood* (1936)—achieve a wider response, sometimes because of their daring in the exploitation of sex or obscenity, more often because of a vitality shared with more demotic fiction. The duplicated typescript or the subsidized periodical, rather than the commercially produced book, is the communication medium for the truly hermetic novel.

The novel that achieves commercial publication but whose limited appeal precludes large financial success can frequently become the object of cult adulation. In the period since World War II, especially in the United States, such cults can have large memberships. The cultists are usually students (who, in an era of mass education, form a sizable percentage of the total population of the United States), or fringes of youth sharing the student ethos, and the novels chosen for cult devotion relate to the social or philosophical needs of the readers. The fairy stories of Tolkien, *The Lord of the Flies* of Golding, the science fiction of Kurt Vonnegut, Jr., have, for a greater or lesser time, satisfied a hunger for myth, symbols, and heterodox ideas, to be replaced with surprising speed by other books. The George Orwell cult among the young was followed by a bitter reaction against Orwell's own alleged reactionary tendencies, and such a violent cycle of adoration and detestation is typical of literary cults. Adult cultists tend, like young ones, to be centred in universities, from which they circulate newsletters on *Finnegans Wake*, Anthony Powell's *Music of Time* sequence, and the works of Evelyn Waugh. Occasionally new public attention becomes focussed on a neglected author through his being chosen as a cult object. This happened when the novellas of Ronald Firbank, the anonymous comic novel *Augustus Carp, Esq.*, and G.V. Desani's *All About Mr. Hatter* got back into print because of the urging of minority devotees. Despite attempts to woo a larger public to read it, Malcolm Lowry's *Under the Volcano* obstinately remained a cult book, while the cultists performed their office of keeping the work alive until such time as popular taste should become sufficiently enlightened to appreciate it.

Detective, mystery, thriller. The terms detective story, mystery, and thriller tend to be employed interchangeably. The detective story thrills the reader with mysterious crimes, usually of a violent nature, and puzzles his reason

Chosisme
and
tropisme

Student
cultists

until their motivation and their perpetrator are, through some triumph of logic, uncovered. The detective story and mystery are in fact synonymous, but the thriller frequently purveys adventurous *frissons* without mysteries, like the spy stories of Ian Fleming, for example, but not like the spy stories of Len Deighton, which have a bracing element of mystery and detection. The detective novel began as a respectable branch of literature with works like Poe's "Murders in the Rue Morgue" (1841), Dickens' unfinished *Edwin Drood* (1870), and Wilkie Collins' *Moonstone* (1868) and *Woman in White* (1860). With the coming of the Sherlock Holmes stories of Sir Arthur Conan Doyle, at the beginning of the 20th century, the form became a kind of infraliterary subspecies, despite the intellectual brilliance of Holmes's detective work and the high literacy of Doyle's writing. Literary men like G.K. Chesterton practiced the form on the margin, and dons read thrillers furtively or composed them pseudonymously (e.g., J.I.M. Stewart, reader in English literature at Oxford, wrote as "Michael Innes"). Even the British poet laureate, C. Day Lewis, subsidized his verse through writing detective novels as "Nicholas Blake." Dorothy L. Sayers, another Oxford scholar, appeared to atone for a highly successful career as a mystery writer by turning to religious drama and the translating of Dante, as well as by making her last mystery novel—*Gaudy Night* (1935)—a highly literary, even pedantic, confection.

Such practitioners as Agatha Christie, Ellery Queen, Erle Stanley Gardner, Raymond Chandler, to say nothing of the highly commercial Edgar Wallace and Mickey Spillane, have given much pleasure and offended only the most exalted literary canons. The fearless and intelligent amateur detective, or private investigator, or police officer has become a typical hero of the modern age. And those qualities that good mystery or thriller writing calls for are not to be despised, since they include economy, skillful sustention of suspense, and very artful plotting.

The mystery novel was superseded in popularity by the novel of espionage, which achieved a large vogue with the James Bond series of Ian Fleming. Something of its spirit, if not its sadism and eroticism, had already appeared in books like John Buchan's *Thirty-nine Steps* and the "entertainments" of Graham Greene, as well as in the admirable novels of intrigue written by Eric Ambler. Fleming had numerous imitators, as well as a more than worthy successor in Len Deighton. The novels of John Le Carré found a wide audience despite their emphasis on the less glamorous, often even squalid aspects of international espionage; his works include *The Spy Who Came in from the Cold* (1963) and *Smiley's People* (1980).

Western. Man's concern with taming wild land, or advancing frontiers, or finding therapy in reversion from the civilized life to the atavistic is well reflected in adventure novels, beginning with James Fenimore Cooper's novels of the American frontier *The Pioneers* (1823) and *The Last of the Mohicans* (1826). As the 19th century advanced, and new tracts of America were opened up, a large body of fiction came out of the men who were involved in pioneering adventure. Mark Twain's *Roughing It* (1872) may be called a frontier classic. Bret Harte wrote shorter fiction, like "The Luck of Roaring Camp" (1868), but helped to spread an interest in frontier writing to Europe, where the cult of what may be termed the western novel is as powerful as in America. Owen Wisters' *Virginia* (1902), Andy Adams' near-documentary *Log of a Cowboy* (1903), Emerson Hough's *Covered Wagon* (1922), from which the first important western film was made in 1923, Hamlin Garland's *Son of the Middle Border* (1917), and O.E. Rølvaag's *Giants in the Earth* (1927) all helped to make the form popular, but it is to Zane Grey—who wrote more than 50 western novels—that lovers of frontier myth have accorded the greatest devotion. The western is now thought of predominantly as a cinematic form, but it arose out of literature. Other frontier fiction has come from another New World, the antipodes—South Africa as well as the Australian outback—but the American West has provided the best mythology, and it is still capable of literary treatment. Sophisticated literary devices may be grafted onto the western—surrealistic fantasy or parallels

to Shakespeare or to the ancient classics—but the peculiar and perennial appeal of the western lies in its ethical simplicity, the frequent violence, the desperate attempt to maintain minimal civilized order, as well as the stark, near-epic figures from true western history, such as Billy the Kid, Calamity Jane, Wyatt Earp, Annie Oakley, and Jesse James.

The best-seller. A distinction should be made between novels whose high sales are an accolade bestowed on literary merit and novels that aim less at aesthetic worth than at profits. The works of Charles Dickens were best-sellers in their day, but good sales continue, testifying to a vitality that was not purely ephemeral. On the other hand, many best-selling novels have a vogue that is destined not to outlast the time when they were produced. It is a characteristic of this kind of best-seller that the writing is less interesting than the content, and that the content itself has a kind of journalistic oversimplification that appeals to unsophisticated minds. The United States is the primary home of the commercial novel whose high sales accrue from careful, and sometimes cold-blooded, planning. A novel in which a topical subject—such as the Mafia, or corruption in government, or the election of a new pope, or a spate of aircraft accidents, or the censorship of an erotic book—is treated with factual thoroughness, garnished with sex, enlivened by quarrels, fights, and marital infidelities, presented in nonliterary prose, and given lavish promotion by its publisher may well become a best-seller. It is also likely to be almost entirely forgotten a year or so after its publication. The factual element in the novel seems to be necessary to make the reader feel that he is being educated as well as diverted. Indeed, the conditions for the highest sales seem to include the reconciliation of the pornographic and the didactic.

A novel with genuine aesthetic vitality often sells more than the most vaunted best-seller, but the sales are more likely to be spread over decades and even centuries rather than mere weeks and months. The author of such a book may, in time, enrich others, but he is unlikely himself to attain the opulence of writers of best-sellers such as Harold Robbins or Irving Wallace.

Fantasy and prophecy. The term science fiction is a loose one, and it is often made to include fantastic and prophetic books that make no reference to the potentialities of science and technology for changing human life. Nevertheless, a novel like Keith Roberts' *Pavane* (1969), which has as a premise the conquest of England by Spain in 1588, and the consequent suppression rather than development of free Protestant intellectual inquiry, is called science fiction, though such terms as "fiction of hypothesis" and "time fantasy" would be more fitting. The imaginative novelist is entitled to remake the existing world or present possible future worlds, and a large corpus of fiction devoted to such speculative visions has been produced in the last hundred years, more of it based on metaphysical hypotheses than on scientific marvels. Jules Verne and H.G. Wells pioneered what may be properly termed science fiction, mainly to an end of diversion. Since the days of Wells's *Time Machine* (1895) and *Invisible Man* (1897), the fiction of hypothesis has frequently had a strong didactic aim, often concerned with opposing the very utopianism that Wells—mainly in his nonfictional works—built on the potentialities of socialism and technology. Aldous Huxley's *Brave New World* (1932) showed how dangerous utopianism could be, since the desire for social stability might condone conditioning techniques that would destroy the fundamental human right to make free choices. Toward the end of his life Huxley produced a cautious utopian vision in *Island* (1962), but the dystopian horrors of his earlier novel and of his *Ape and Essence* (1948) remain more convincing. Orwell's *Nineteen Eighty-four* (1949) showed a world in which a tyrannic unity is imposed by a collective solipsism, and contradictions are liquidated through the constant revision of history that the controlling party decrees. Anthony Burgess' *Clockwork Orange* (1962) and *Wanted Seed* (1962) portray ghastly futures that extrapolate, respectively, philosophies of crime control and population control out of present-day tendencies that are only potentially dangerous.

The enduring and the ephemeral

A large number of writers practice prophetic fantasy with considerable literary skill and careful factual preparation—Kurt Vonnegut, Jr., Ray Bradbury, Italo Calvino, Isaac Asimov, J.G. Ballard, to name only a few—and novelists whose distinction lies mainly in more traditional fields have attempted the occasional piece of future-fiction, as in the case of L.P. Hartley with his *Facial Justice* (1961) and Evelyn Waugh in *Love Among the Ruins* (1953). The fantasist who fantasizes without prophetic or warning intent is rarer, but works such as Nabokov's *Invitation to a Beheading*, Tolkien's *Lord of the Rings* cycle, and Christine Brooke-Rose's *Out* (1964) represent legitimate and heartening stretching of the imagination, assurances that the novelist has the right to create worlds, as well as characters, of his own. However, the dystopian novel can have a salutary influence on society, actively correcting regressive or illiberal tendencies, and *Brave New World* and *Nineteen Eighty-four* can be cherished as great didactic landmarks, not just as works of literary art.

Proletarian. The novel that, like Dickens' *Hard Times* (1854), presents the lives of workingmen or other members of the lower orders is not necessarily an example of proletarian fiction. The category properly springs out of direct experience of proletarian life and is not available to writers whose background is bourgeois or aristocratic. Consequently, William Godwin's *Caleb Williams* (1774), or Robert Bage's *Hermesprung* (1796), although, like *Hard Times*, sympathetic to the lot of the oppressed worker, is more concerned with the imposition of reform from above than with revolution from within, and the proletarian novel is essentially an intended device of revolution. The Russian Maksim Gorky, with works like *Foma Gordeyev* (1900) and *Mother* (1907), as well as numerous short stories portraying the bitterness of poverty and unemployment (in fact, the pseudonym *Gorky* means "Bitter"), may be taken as an exemplary proletarian writer. The United States has produced a rich crop of working-class fiction. Such Socialist writers as Jack London, Upton Sinclair, John Dos Passos, and Edward Dahlberg, however, did not witness the triumph of the workers' revolution in their own country, as Gorky did in his, and it is the fate of the American proletarian novelist, through literary success, either to join the class he once dreamed of overthrowing or to become anarchic and frustrated. In the Soviet Union the proletarian novel was doomed to disappear in the form that Gorky knew, for it is the essence of the revolutionary novel to possess vitality and validity only when written under capitalist "tyranny."

England has produced its share of working-class novelists exuding bitterness, such as Alan Sillitoe, with his *Saturday Night and Sunday Morning* (1958), but conditions apt for revolution have not existed in Britain for more than a century. British novelists who emerged after World War II, such as John Braine (*Room at the Top*), Keith Waterhouse (*There Is a Happy Land*), Kingsley Amis (*Lucky Jim*), and Stan Barstow (*A Kind of Loving*), provided a solution to working-class frustration in a fluid system of class promotion: revolution is an inadmissible dream. Generally speaking, in the novel, which is preoccupied with individuals rather than with groups, it is difficult to make the generalized political statements that are meat and drink to the revolutionary propagandist.

Other types. The categories briefly discussed above are among the most common fictional forms. Theoretically there is no limit to the number available, since changing social patterns provide fresh subjects and fresh taxonomies, and new metaphysical and psychological doctrines may beget new fictional approaches to both content and technique.

Other categories of fictional art include the erotic novel (which may or may not be pornographic), the satirical novel, the farcical novel, the novel for or about children, the theological novel, the allegorical novel, and so on. Types of fiction no longer practiced, since their real-life referents no longer exist, include the colonial novel—such as E.M. Forster's *Passage to India* (1924), Henri Fauconnier's *Malaisie* (1930), and the African sequence of Joyce Cary—and space fantasy like H.G. Wells's *First Men in the Moon* (1901). One may read examples of a departed

category with pleasure and profit, but the category can no longer yield more than parody or pastiche.

New kinds of fiction fill in the gaps, like the novel of negritude, the structuralist novel (following the linguistic sociologists and anthropologists), the homosexual novel, the novel of drug hallucination, and so on. So long as human society continues to exist, the novel will exist as its mirror, an infinitude of artistic images reflecting an infinitude of life patterns.

THE NOVEL IN ENGLISH

The United Kingdom. England's chief literary achievements lie in the fields of drama and poetry, and the attitude of English novelists to their form was, for a long time, cheerfully empirical and even amateurish. Elizabethan novels, *novelle* rather, imitated the Spanish picaresque story, and Thomas Nashe's *Unfortunate Traveller* (1594) is a good, bustling, vital example of a rapidly composed commercial work more concerned with sensational incident and language than with shape or character. Daniel Defoe (1660?–1731) is often considered to be the true progenitor of the long English novel, but his *Robinson Crusoe* and *Moll Flanders* are loosely constructed, highly episodic, and presented as mock biography rather than real fiction. It is with *Pamela*, by Samuel Richardson (1689–1761), that the tradition of serious, moral fiction in English may be said to begin, but later 18th-century novelists reverted to the picaresque and comic. Henry Fielding (1707–54) wrote his first novel, *Joseph Andrews*, as parody on *Pamela*, but his masterpiece, *Tom Jones*, is original if shapeless, an example of English literary genius sinking thankfully back into the casual and improvisatory. Laurence Sterne (1713–68) produced a great mad work, *Tristram Shandy*, that, in its refusal to truckle to any rules of structure, remains still a quarry for avant-garde novelists; and Tobias Smollett (1721–71) wrote picaresque satire in *Roderick Random* and *Peregrine Pickle*, full-blooded portraits of the age that impress more with their vigour than with their art.

The Romantic Age brought, rather paradoxically, the cool and classically shaped novels of Jane Austen (1775–1817), a major practitioner and still a model for apprentices in the craft. Sir Walter Scott (1771–1832), a Scotsman who wrote about romantic historical Scotland, must be regarded as an international figure whose influence was greater even than Richardson's, since he more than anyone established the historical novel as the primary fictional form in Europe. In Scott's work, nevertheless, the traditional faults of the British novel may be desecrated—namely, episodic formlessness, an ebullience of texture rather than a clean narrative line.

The Victorian Age moved out of the romantic past, or, as with William Makepeace Thackeray (1811–63), stayed with it only to deromanticize it. Charles Dickens (1812–70) was indebted to the picaresque tradition but turned reformist eyes on his own age. Thackeray and Dickens are complementary in that the first attacks the upper classes while the second showers sympathy, sometimes of a mawkish kind, on the lower classes. With George Eliot (1819–80), the first true English psychological novels appear, strong in their moral content, and George Meredith (1828–1909) may be said to have anticipated—in *The Ordeal of Richard Feverel* and *The Egoist*—the approach in depth that characterized the psychological novel of the 20th century.

Both Charlotte Brontë (1816–55) and Emily Brontë (1818–48) exemplify the capacity of the English novel to achieve solitary "sports" unrelated to any current or tradition. Both *Wuthering Heights*—a superb evocation of the soul of a locality, with a love story that is fierce and primitive but recounted with poetic sophistication—and *Jane Eyre*, an exceedingly frank and still shocking study of a love that rides over Victorian conventions, are unlike any other books of their time, or of any other time, though their qualities have been diluted into hundreds of popular 20th-century romances. Later Victorians, particularly Samuel Butler (1835–1902) and Thomas Hardy (1840–1928), reflected those changes in the educated English sensibility that had been brought about by the new science.

Reform
versus
revolution

The
Romantic
and
Victorian
ages

Outmoded
types

Hardy's world is one in which the Christian God has been replaced by a malevolent Providence—the poet-novelist's theologization of scientific determinism. Butler's *Way of All Flesh*, a work that contrives to be both bitterly realistic and highly comic, demonstrates the working of Darwinian evolution in social institutions such as the family and even the church. In many ways, Butler led English fiction into the modern age.

John Galsworthy (1867–1933) showed himself interested in the processes by which old institutions, such as a great Whig family, decay as history advances, but in both style and characterization he is firmly set in the Victorian Age. Of Arnold Bennett (1867–1931) it may be said that he brought to a kind of Thackerayan social realism something of the spirit of the French novel, particularly the anglicized tones of Balzac and Zola. W. Somerset Maugham (1874–1965), though his most considerable achievements lie in the short-story form, wrote novels, like *Of Human Bondage*, that are infused with the delicious acerbities of French naturalism, and all his work, long or short, is exalted by the influence of Guy de Maupassant. Early 20th-century British fiction needed the impact of an alien tradition to jolt it out of bourgeois empiricism, and perhaps the two major influences were both foreigners who had elected to write fiction in England—Joseph Conrad (1857–1924), Polish-born, to whom English was a second language, and Henry James (1843–1916), an American who had drunk deeply at French fountains and brought to the exercise of his craft a scrupulousness and a concern with aesthetic values that was almost obsessive and, it may be said, very un-English. James's influence on such major English-born novelists as Virginia Woolf (1882–1941), E.M. Forster (1879–1970), and Ford Madox Ford (1873–1939) was in the direction of that concern with style which the English novel, unlike the French, has always tried to resist, and this influence remains a potent one on succeeding generations of novelists, not only in England.

Other important Englishmen have remained more interested in content than in style, though in Graham Greene, a Catholic convert, who is primarily known for his "theological" content, there is a very interesting attempt to bring to this a Conradian concern for the solitary-man theme and a Jamesian preoccupation with style. D.H. Lawrence (1885–1930) remains the great modern English novelist who reconciles a highly traditional style (in *The Rainbow*, for example, he often resembles George Eliot) with a subject matter that is revolutionary in the profundity of its human relationships. And Aldous Huxley (1894–1963), though he indulged in formal experiment in *Point Counter Point* and *Eyeless in Gaza*, was always essentially a didactic writer who used the novel form somewhat casually. The same may be said of George Orwell (1903–50), who eschewed stylistic complexity in the interest of a clear message. On the other hand, Evelyn Waugh (1902–66) and Anthony Powell, in a cautious and conservative way, consulted the claims of allusive and evocative prose, as did Wyndham Lewis (1884–1957), who also brought to his novels the aesthetic of an alien art—that of painting.

Novelists like Kingsley Amis and Angus Wilson are seen to belong to traditions already old in the time of Samuel Butler. Amis derives from Fielding (Amis' *Take a Girl Like You* has a moral quality reminiscent of Fielding's *Amelia*), and Wilson never pretends to be other than a disciple of novelists like George Eliot and Dickens. The Victorian kind of novel, as practiced by writers like J.B. Priestley, satisfies large numbers of British readers, and Jamesian scrupulousness and Joycean experiment alike are regarded with amiable suspicion. Such Englishmen as have experimented in fictional technique have usually found a larger audience abroad than at home. Thus Lawrence Durrell found his *Alexandria Quartet* (1962) hailed as a masterpiece in France, while English readers merely liked or disliked it. *A Clockwork Orange* (1962) by Anthony Burgess achieved a large readership in the United States, but it recorded little positive response in the country of its origin.

The fiction of the British Isles is not clearly differentiated into national or regional groups, and to speak of Irish novelists is often to do no more than refer to an irrelevant

place of birth. Oscar Wilde (1854–1900) was more French than Irish; his one novel, *The Picture of Dorian Gray*, has no progenitors in English fiction, but it owes much to the late 19th-century French novelist Joris-Karl Huysmans. George Moore (1852–1933) was born in County Mayo, Ireland, but studied art in Paris, and works like *A Modern Lover* and *A Mummer's Wife* were intended to teach the British novel how to absorb the spirit of Flaubert and Zola. Moore's *Lake*, written under the influence of the Irish literary revival, is dutifully set in Ireland, but this essentially international writer was quick to move back to France for *Héloïse and Abélard* and on to Palestine for *The Brook Kerith*. James Joyce (1882–1941), conceivably the greatest novelist in English in the 20th century, never moves from his native Dublin in any of his fiction, but no less parochial writer can well be imagined. *Ulysses*, in the depth of its historical coverage, may be regarded as one of the last great artistic monuments of the Austro-Hungarian Empire, and *Finnegans Wake* was stimulated by the avant-garde climate of Paris, the resting place of so many Irish "wild geese." It is only perhaps negatively that Ireland shaped Joyce's literary personality; the oppressiveness of the new nationalism made him react in the direction of internationalism, and it eventually forced him to a life-long exile. Flann O'Brien (1911–66), possibly Joyce's true successor, stayed in Dublin and drew, in works like *At Swim-Two-Birds* (1939) and *The Hard Life* (1961), on the native demotic experience, but his techniques came from Europe, not from the Anglo-Irish bourgeois stockpot.

The Scottish novel hardly exists as a national entity; there is nothing in fiction that matches the exploitation of Lallans in the poetry of Burns or of Hugh MacDiarmid (C.M. Grieve). Scott, as has been indicated, belongs to European literature, and later novelists like Robert Louis Stevenson (1850–94) and Sir James Barrie (1860–1937), though Scottish themes and settings are featured triumphantly in their fiction, are essentially men whose literary metropolis was London rather than Edinburgh.

Wales sedulously cultivates Welsh as a living artistic medium, but few novels by Welshmen reach the English or international market, unless—like Emyr Humphreys—they write their works in English as well as Welsh. The English-language novel from Wales hardly exists, except in the form of best-sellers like Richard Llewellyn's *How Green Was My Valley* (1940)—a deliberate capitalization on Welsh picturesqueness—and the exceedingly Joycean prose of the poet Dylan Thomas (1914–53), whose one novel, *Adventures in the Skin Trade*, has very little of Wales in it.

The United States. If the American novel appears to begin with the *Wieland* and *Edgar Huntly* of Charles Brockden Brown (1771–1810), then the fiction of the mother country may be said to have no start on that of the newly independent daughter. Admittedly, Brown owes something to the English Gothic novel, but already a typically American note is struck in the choice of a violent and bizarre form that, in various mutations, was to prove fruitful in the development of American fiction. James Fenimore Cooper (1789–1851) struck out, in the pentateuch called the "Leatherstocking" tales, in another direction suitable to the American genius and experience: *The Last of the Mohicans* and *The Prairie*, though their prose style is perhaps overelaborate, are full of the spirit of a young nation confronting the wilderness and advancing its frontiers. Harriet Beecher Stowe (1811–96) produced, in *Uncle Tom's Cabin*, an antislavery novel whose sensational devices owe something to the Gothic movement but whose style is rooted in that European romantic tradition fathered by Scott. Even Edgar Allan Poe (1809–49) and Nathaniel Hawthorne (1804–64) are incompletely emancipated from Europe, though Hawthorne's *Scarlet Letter* shows a preoccupation with sin that finds no fictional counterpart in the mainstream of the European novel. America, aware of the darkness and mystery of a land still mainly undiscovered, was correspondingly aware, in both Gothic and eschatological fiction, of the dark places of the mind.

The true emancipation from Europe comes in the works of Hawthorne's friend Herman Melville (1819–91), whose

Irish and
Scottish
novels

Impact of
non-British
writers

Moby Dick creates a totally American fusion by combining plain adventure with a kind of Manichaean symbolism—implying a belief that the universe is under the dominion of two opposing principles, one good and the other evil. Mark Twain (1835–1910) brought the Mississippi frontier region into the literary geography of the world with *Tom Sawyer* and *Huckleberry Finn*, combining very American humour with harsh social criticism, as also in the unique *Pudd'nhead Wilson*. From Twain on, the new regions of the United States become the materials of a wealth of fiction. Thus, while Bret Harte (1836–1902) wrote about California and George W. Cable (1844–1925) on Louisiana, Indiana was celebrated in *The Hoosier Schoolmaster* of Edward Eggleston (1837–1902) and the *Penrod* stories of Booth Tarkington (1869–1947).

But America could not wholly turn its back on Europe, and the early 20th century was characterized by the discovery of the phenomenon called American "innocence" in confrontation with the decadent wisdom and sophistication of the Old World. If, in *Innocents Abroad*, Mark Twain and his fellow travellers remained unimpressed by Europe, Henry James (1843–1916) devoted millions of words to the response of impressionable Americans to the subtle deciduous culture their ancestors had abandoned in the search for a new life. Zolaesque realism entered the American novel with William Dean Howells (1837–1920). Theodore Dreiser (1871–1945) shocked his audiences with the candour and pessimism of *Sister Carrie* and *An American Tragedy*, while Stephen Crane (1871–1900) and Frank Norris (1870–1902) made realism enter, respectively, the traditionally romantic field of war (*The Red Badge of Courage*) and the pastoral life of the promised land of California (*The Octopus*). American innocence ceased to be a fictional property. The brutality in Jack London (1876–1916) joins hands with the depictions of "natural man" in the works of Ernest Hemingway (1899–1961), while the urban fiction of James T. Farrell (1904–79) and John O'Hara (1905–70) has a wholly American realism that leaves the naturalism of Zola far behind.

The 20th-century American novel is "compartmentalized" in a manner that seems to give the lie to any allegation of cultural unity. The cult of the regional novel has continued. Ohio came under the microscopic scrutiny of Sherwood Anderson (1876–1941) in *Winesburg, Ohio*, just as the small-minded hypocrisy and materialism of the Midwest fired the brilliant sequence of Sinclair Lewis (1885–1951) from *Main Street* on. But the best regional fiction has come from the South, with William Faulkner (1897–1962), whose technical master is Joyce, Ellen Glasgow (1874–1945), Flannery O'Connor (1922–64), Eudora Welty, Erskine Caldwell, and Robert Penn Warren. The immigrant communities of Nebraska were the subject of *My Antonia*, by Willa Cather (1873–1947), while O.E. Rølvaag (1876–1931) dealt with the South Dakota Norwegians in *Giants in the Earth*. The urban Jew has become the very spokesman of the contemporary American experience with writers like Saul Bellow, Herbert Gold and Bernard Malamud; while the black American has been made highly articulate in works like *Go Tell It on the Mountain* by James Baldwin, *Invisible Man* by Ralph Ellison, and *Home to Harlem* by Jamaican-born Claude McKay (1890–1948).

Meanwhile, the fiction of social criticism that derives from Howells and was made popular by men like Upton Sinclair (1878–1968) goes on, either in "muckraking" best-sellers or in such specialized attacks on the American ethos as are exemplified by war books like Joseph Heller's *Catch-22* and Kurt Vonnegut's *Slaughterhouse-Five*. American realism has gained a reputation for candour unequalled in any other literature of the world, so that the *Tropic of Cancer* and *Tropic of Capricorn* of Henry Miller (1891–1980) remain models for sexual frankness that encouraged writers like Hubert Selby, Jr. (*Last Exit to Brooklyn*), and Philip Roth (*Portnoy's Complaint*) to uncover areas of eroticism still closed to the European novel.

To balance the concern with naked subject matter, American fiction has also fulfilled the promise of its earlier expatriates beginning with Gertrude Stein (1874–1946) to match the French in a preoccupation with style. The

distinction of F. Scott Fitzgerald (1896–1940) lay in the lyrical intensity of his prose as much as in his "jazz age" subject matter, and writers like Truman Capote and John Updike have dedicated themselves similarly to perfecting a prose instrument whose effects (like those of Joyce) touch the borders of poetry. Perhaps the glories and potentialities of American fiction are best summed up in the novels of Vladimir Nabokov (1899–1977). His early works belonged to Europe, but when he took to writing in English he created a sophisticated, cosmopolitan, and highly poetic style. His work, as in *Lolita* and *Pale Fire*, concentrates with unparalleled intensity on the immediacies of American life in the 20th century.

The British Commonwealth. Canada has two literatures—one in French as well as one in English—and, in the taxonomy of this article, it will be as well to forget the course of history and consider French-Canadian novelists along with those of the separated, and officially abandoned, mother country. The Canadian novel in English begins, appropriately enough, with a Richardson—John Richardson (1796–1852), author of *Wacousta*, a story of the Indian uprising led by Pontiac. But, with the exception of James DeMille (1836–80), author of a remarkable novel, *A Strange Manuscript Found in a Copper Cylinder*, and William Kirby (1817–1906), whose *Golden Dog* is an interesting long romance of 18th-century Quebec, no novelist of stature emerged in what was a great period of literary expansion in the United States. The historical novels of Sir Gilbert Parker (1862–1932), the western romances of Ralph Connor (1860–1937), and the world-famous *Anne of Green Gables* by Lucy Maud Montgomery (1874–1942) exhibit a stylistic tameness, along with that jejune oversimplification of psychology that is a mark of the deliberately "popular" novel.

More distinctive Canadian work came in the early 20th century, as in the prairie novels of Robert Stead (1880–1959) and Frederick Philip Grove (1871–1948). Morley Callaghan, a chronicler of urban life, may be seen as a writer approaching international stature, and that tough "European" realism that braces his work is also to be found, along with a rather exotic elegance, in the novels of Robertson Davies. Although some younger novelists, like Margaret Atwood and Alice Munro, deal with growing up and rites of passage in Canada, others, like Brian Moore and Mordecai Richler, have sought artistic stimulation outside Canada, and they show general uneasiness about the lack of a cultivated audience in their great but underpopulated land.

A similar complaint may be heard in Australia and New Zealand, where the fictional tradition has long accommodated itself to the simple literary needs of unsophisticated settlers. The first true Australian novel was probably *For the Term of His Natural Life*, by Marcus Clarke (1846–81), which dealt, appropriately, considering the penal origins of the Australian settlements, with life in prison; and Rolf Boldrewood, the pseudonym of Thomas A. Browne (1826–1915), had a cognate approach of excessive simplicity and melodrama to the days of the gold rush in his *Robbery Under Arms*. Australian farm life was the theme of the novels that Arthur Hoey Davis (1868–1935) wrote under the pen name of Steele Rudd, beginning with *On Our Selection*. But something that can only be termed native Australianism—which has less to do with trades and scenic backgrounds than with a whole new language and collective philosophy of life—appears for the first time in *My Brilliant Career* (1901), by Miles Franklin (1879–1954). The expression of a national personality is not, however, enough.

Progressive 20th-century Australian novelists look hungrily toward Europe, aware that native conservatism will not support the kind of experimentalism in literature that it is prepared to take for granted in painting and architecture. The most considerable modern Australian novelist is Patrick White, but even such massive achievements as *Voss* and *Riders in the Chariot* are set firmly in the Dostoyevskian tradition. Morris West is Australian only by birth; he seeks international themes and the lure of the American best-seller market. It seems that an Australian writer can succeed only if he renounces his Australianism

American
innocence
and
European
sophistica-
tion

Australian
and New
Zealand
novels

and goes into exile. It is still a source of humiliation to Australian literateurs that the finest novel about Australia was written by a mere visitor—*Kangaroo* (1923), by D.H. Lawrence.

New Zealand, with its much smaller population, presents even greater problems for the native novelist seeking an audience, and this—along with New Zealand's closer tie with the mother country—explains why such fiction writers as Katherine Mansfield (1888–1923) and Dame Ngaio Marsh (1899–1982) made London their centre. Young novelists like Janet Frame, Ian Cross, and Maurice Shadbolt show, however, a heartening willingness to render the New Zealand scene in idioms and techniques more progressive than anything to be found in Australia, and they are prepared to resist the temptations of the expatriate life.

Inevitably, territories like India and Africa entered the federation of English literature very much later than those dominions founded on the English language. With few exceptions, African and Indian novelists employ English as a second language, and one of the charms of their novels lies in a creative tension between the adopted language and the native vernacular (needless to say, this is usually self-consciously exploited—often for poetic, but more frequently for comic, effect).

Of Indian novelists, R.K. Narayan, in works like *The English Teacher* (U.S. title *Grateful to Life and Death*) and *The Man-Eater of Malgudi*, exhibits an individual combination of tenderness and humour, as well as a sharp eye for Indian foibles. Raja Rao, whose best known novel is *The Serpent and the Rope*, achieves remarkable prosodic effects through allowing Sanskrit rhythm and idiom to fertilize English. Of younger Indian novelists, Balachandra Rajan is notable, in *The Dark Dancer* and *Too Long in the West*, for an ability to satirize the Anglo-American way of life with the same suave elegance that informs his tragicomic view of the East. Khushwant Singh presents, in *I Shall Not Hear the Nightingale*, a powerful chronicle of Sikh life during that period of imperial dissolution that began with World War II. English seems established as the medium for the Indian novel, and it is interesting to note an ability on the part of nonnative Indian residents who have practiced the form to be absorbed into a new "Indo-English" tradition. Rudyard Kipling's *Kim* is respected by Indian writers, and E.M. Forster's *Passage to India* is a progenitor of one kind of Indian novel. The novels of Paul Scott (1920–78)—such as the tetralogy *The Raj Quartet* (1976)—spring out of a love of the country and an understanding of its complexities not uncommon among former British soldiers and administrators. The Indian novel is perhaps a product of territory rather than of blood.

The most important of the new African writers come from the West Coast, traditional fount of artists, and they are mostly characterized by immense stylistic vigour, a powerful realism, and, often, a satirical candour unsoftened by the claims of the new nationalism. Chinua Achebe, in *Things Fall Apart* and *No Longer at Ease*, renders remarkably the tones of Umuaro speech and thought, and exhibits, as also in *Arrow of God*, a concern for that rich native culture whose extirpation is threatened by imported Western patterns of life and government. His *Man of the People* deals sharply with corruption and personality cult in a newly independent African state modelled on his own Nigeria. A fellow Nigerian, Amos Tutuola, has gained an international reputation with *The Palm-Wine Drinkard*—a richly humorous novel permeated with the spirit of folklore. Cyprian Ekwensi is best known for his *Jagua Nana*, a wry study of the impact of the new materialism (symbolized by the "Jagua" ear of the title) on the tribal mind. Onuora Nzekwu, in *Blade Among the Boys*, has a graver theme: the conflict between Ibo religion and imported Christianity in the upbringing of a sensitive and confused young man.

South African novels have, traditionally, dealt with those pioneering themes still exemplified in the work of Laurens van der Post and Stuart Cloete, but the territory has made its entry into world literature comparatively recently chiefly because the official racist policies have inspired a powerful fiction of protest. William Plomer (1903–73) may be regarded as the father of anti-apartheid literature,

although his *Turbott Wolfe* appeared in 1925, long before the state doctrine was articulated. The theme of this novel was the necessity for white and black blood to mix and ensure a liberal South African future not given over purely to white domination. Alan Paton, in *Cry, the Beloved Country*, has produced the most popular novel of protest, but Dan Jacobson, Nadine Gordimer, and especially Doris Lessing have amplified mere protest into what may be termed a kind of philosophical fiction, often distinguished enough to rank with the best work of Europe and America. Lessing's sequence *Children of Violence* finds in the South African system of government a starting point for denunciation of wrongs that turn out to be social and sexual as well as racial.

The varied fictional achievements of the Caribbean are large. Trinidad has produced the two best known West Indian writers—Samuel Selvon and V.S. Naipaul, both of East Indian descent. Selvon's work (*A Brighter Sun, An Island Is a World, The Lonely Londoners*) is grim, bitter, capable of vivid evocation of the Trinidadian scene, but Naipaul, after *A House for Mr. Biswas*, has shown signs of habituation to his English exile, so that his *Mr. Stone and the Knights Companion* seems to be a novel totally nourished by the London world in which it is set. Most Caribbean novelists, finding their publishers and their audiences in England, transfer themselves thither and cut themselves off from all but remembered roots. This is true of Edgar Mittelholzer (1909–65), an ebullient and prolific writer whose later books all had English settings. Wilson Harris has, in English exile, created an astonishing Guianan tetralogy in which poetry and myth and symbolic difficulty have a place. George Lamming and John Hearne are both notable for firm and economical prose and masterly scene painting. The aesthetic prospects for the West Indian novel seem excellent, but the absorption of its practitioners into the larger English-speaking world represents a symptom that plagues the practice of literature in so many Commonwealth countries—the lack of adequate publishing facilities and, more than that, the failure of a cultivated readership to emerge. It is undoubtedly unhealthy for an author to have to seek primary communication with foreigners.

EUROPE

Russian. The Russian novel properly begins with Nikolay Karamzin (1766–1826), who introduced into Russian literature not only the exotic sentimental romanticism best seen in his novels *Poor Lisa* and *Natalya, the Boyar's Daughter* but that large Gallie vocabulary that was to remain a feature of the literary language. But those qualities that best distinguish Russian fiction—critical realism and spirituality—first appeared in the work of Mikhail Lermontov (1814–41), whose *Hero of Our Time* is the pioneering Russian psychological novel. Nikolay Gogol (1809–52), satirizing provincial mores in *Dead Souls*, ushered in, perhaps unintentionally, a whole fictional movement—the "literature of accusation." Vissarion Belinsky (1811–48), the father of Russian literary criticism, formulated the theory of literature in the service of society, and Ivan Turgenev (1818–83), produced a classic "accusatory" novel in *A Sportsman's Sketches*. But he, and the other major novelists who emerged in the mid-19th century, can hardly be considered in terms of literary movements. Fyodor Dostoyevsky (1821–81), with *The Possessed, The Idiot, Crime and Punishment*, and *The Brothers Karamazov*, affirmed with idiosyncratic power the great spiritual realities, and Leo Tolstoy (1828–1910) produced two of the greatest novels of all time—*War and Peace* and *Anna Karenina*, revelatory of the Russian soul but also of the very nature of universal man and human society.

The later days of the 19th century saw a shift in fictional radicalism. The theories of Karl Marx influenced the accusatory writers in the direction of the plight of the proletariat, not, as had been the old way, that of the peasantry. Maksim Gorky (1868–1936), Leonid Andreyev (1871–1919), Aleksandr Kuprin (1870–1938), and the Nobel Prize winner (1933) Ivan Bunin (1870–1953) wrote of the Russian urban experience and helped to create the literary climate of the Soviet regime. Generally, since

English as a second language of the novelist

African novelists

Caribbean works

The influence of the Soviet government

the beginning of the first five-year plan in 1928, there had been a division between what the regime regarded as valuable in the practice of the novel and what the rest of the world thought. Mikhail Sholokhov (1905–84) depicted, with no evident propagandist slanting, the Revolution and civil war in *Quiet Flows the Don*, and Fyodor Gladkov (1883–1958) was one of the few novelists readable outside the ranks of the Soviet devout in the new category of economic or industrial fiction. The Russian novel of the age following World War II, however, underwent a dramatic schism, in which writers like Aleksandr Solzhenitsyn and Boris Pasternak (1890–1960) could receive the Nobel Prize but be officially condemned in the Soviet Union, while the fictional darlings of the regime were recognized, in the non-Communist world, as possessing little or no aesthetic merit.

German. Goethe (1749–1832), who practiced so many arts with such notable brilliance, may be regarded as the first major novelist of Germany. His life covers the whole period of the Enlightenment, with its insistence on a national literary spirit, the *Sturm und Drang* movement, and that phase of *Weltschmerz* which *The Sorrows of Young Werther* fostered. It was *Werther* more than any other work that carried German literature into the world arena; it remained influential in Europe when German Romanticism had already burned itself out. The German reaction against Romanticism was expressed in the regionalism of Theodor Storm (1817–88) and Fritz Reuter (1810–74), who sought to render with objective fidelity the life of their native provinces (northwestern and northeastern Germany, respectively). Swiss writers like Gottfried Keller (1819–90) belonged to the movement of mainstream German regional realism, featuring the picturesque solidities of Switzerland. But the post-Goethean major achievements in the novel had to wait for the Impressionist movement, which produced the works of Thomas Mann (1875–1955) and Hermann Hesse (1877–1962)—fiction concerned less with a roughly hacked slice of life than with form and aesthetic delicacy. World War I brought Expressionism—the nightmares of the German-Czech Franz Kafka (1883–1924) and the psychoanalytical novels of Jakob Wassermann (1873–1934), with their pleas for humanity and justice.

When the true historical nightmare of the Nazi regime followed—predicted, in a sense, by Kafka and Mann—liberal German fiction was suppressed, and liberal German novelists like Mann and Erich Maria Remarque (1898–1970), author of *All Quiet on the Western Front*, went into exile. The brutal, philistine, and nationalistic novels of the true Nazi novelists exist only as curious and frightening relics of an era of infamy. The task of postwar novelists like Günter Grass—author of *The Tin Drum* and *Dog Years*—and Uwe Johnson has been to diagnose the long sickness and force German fiction into new directions—often with the help of surrealism, irony, and verbal experiment.

French. Although the mid-16th-century satirical piece *Gargantua and Pantagruel* of François Rabelais has had, and is still having, a profound influence on the world novel, it would be wrong to place it in the line of true French fiction, which has always manifested a preoccupation with form, order, and economy—qualities totally un-Rabelaisian. The first notable French novel—*The Princess of Clèves* by Madame de La Fayette (1634–93)—shows none of the vices of the English novels of a century later: it is firmly constructed and takes character seriously, as does *Manon Lescaut* by l'Abbé Prévost (1697–1763). With such works the psychological novel was established in Europe. A keen hold on reality and a concern with the problems of man as a social being animates even the fiction of the Romantic school—works like *Indiana* and *Lélia* by George Sand (1804–76) and the heavyweight romances of Victor Hugo (1802–85). But the true glories of French fiction come with the reaction to romanticism and sentimentality, as exemplified in the novels of Jean-Jacques Rousseau (1712–78). The great fathers of realism are Stendhal (1783–1842), Balzac (1799–1850), and Flaubert (1821–80), and their influence is still active. Émile Zola (1840–1902) moved away from the artistic detachment that Flaubert preached and practiced, but, in his *Chroni-*

cles of the Rougon Macquart Family, he tried to emulate the encyclopaedic approach to the novel of Balzac, whose *Human Comedy* is meant to be a history of society in a hundred episodes. The leader of the Naturalist school, Zola saw human character as a product of heredity and environment.

It was left to the 20th-century French novel to cast doubt on a mechanistic or deterministic view of man, to affirm the irrational element in his makeup, and to emphasize the primacy of the will. The temporal treadmill of Balzac and Zola has no place in the masterpiece of Marcel Proust (1871–1922); *Remembrance of Things Past*, if it has a philosophy, says more about the creative élan vital of Bergson, the human essences that underlie the shifting phenomena of time and space, than the social jungle the realists had taken for reality. André Gide (1869–1951) seems to make a plea for human aloofness from environment so that the essentially human capacity for change and growth may operate. André Malraux (1901–76), the forerunner of the Existentialist novelists, demonstrated, in *Man's Fate*, the necessity for human involvement in action as the only answer to the absurdity of his position in a huge and indifferent or malevolent universe. Jean-Paul Sartre (1905–80) and Albert Camus (1913–60) similarly emphasized man's freedom to choose, to say no to evil, to define himself through action.

Other French novelists have been more concerned with recording the minutiae of human life as a predominantly sensuous and emotional experience, like Colette (1873–1954), or with taking the religious sensibility as a fictional theme, like François Mauriac (1885–1970). The practitioners of the *anti-roman*—Butor, Robbe-Grillet, Nathalie Sarraute—pursue their attempts at liquidating human character in the traditional novelistic sense. Samuel Beckett, an Irishman who has turned himself into a major French stylist, goes his own way, presenting—in works like *Molly*, *Malone Dies*, *Watt*, and *The Unnamable*—mankind reduced to degradation and absurdity but somehow admirable because it survives.

French-Canadian fiction inevitably suffers from comparison with the glories of the mother country. Before 1900 there is little of value to record—except perhaps the historical romance of Philippe de Gaspé (1786–1871), *Les Anciens Canadiens*—but Louis Hémon (1880–1913) produced a genuine classic in *Maria Chapdelaine*, a story of Canadian pioneer life, original, moving, and sensitive. The somewhat provincial character of French-Canadian life, dominated by the Church and by outmoded notions of morality, has not been conducive either to fictional candour or to formal experiment, and metropolitan France is unimpressed, for the most part, by the literature of the separated brethren. But there is time for an avant-garde to develop and great masters to appear.

Spanish. The great age known as “El Siglo de Oro,” or the Golden Age, produced what is conceivably (it must contest this claim with *War and Peace*) the most magnificent of all novels—the *Don Quixote* of Miguel de Cervantes Saavedra (1547–1616). A satire on chivalry that ends as a humane affirmation of the chivalric principle, it encloses—in tender or comic distortion—other fictional forms that flourished in the Spanish Golden Age. The *novela picaresca* fathered a whole European movement, and its best monuments are perhaps the anonymous *Lazarillo de Tormes* (1554), *Guzmán de Alfarache* by Mateo Alemán (1547–c. 1614), and *El diablo cojuelo* (*The Limping Devil*) by Luis Vélez de Guevara (1579–1644). The pastoral novel, another popular but highly stylized form, was less true fiction than a sort of prose poem, in which lovers in shepherd's disguise bewailed an unattainable or treacherous mistress. But here, since the fictional lovers were often real-life personages in the cloak of a rustic name, the germs of the *roman à clef* are seen stirring. The *novela morisca* was an inimitable Spanish form, a kind of fictional documentary about the wars between Christians and Muslims, as in *Guerras civiles de Granada* (*Civil Wars of Granada*) by Ginés Pérez de Hita (c. 1544–c. 1619).

The decline of Spain as a European power is associated, in literature, with feeble nostalgia for the Golden Age or feebler imitations of French classicism. The Spanish novel

French-Canadian novels

The Golden Age in Spain

Suppression of the novel by the Nazis

began to reemerge only in the early 19th century, when a journalism celebrating regional customs encouraged the development of the realistic regional novel, with Fernán Caballero (1796–1877), Armando Palacio Valdés (1853–1938), and the important Vicente Blasco Ibáñez (1867–1928), whose *Sangre y arena* (*Blood and Sand*), *Mare nostrum*, and *Los cuatro jinetes del Apocalipsis* (*The Four Horsemen of the Apocalypse*) achieved universal fame and were adapted for the screen in the 1920s.

Generation
of '98

The fiction of the Generation of '98—which took its name from the year of the Spanish-American War, a cataclysmic event for Spain that bereaved it of the last parts of a once great empire—wasted no time on national nostalgia or self-pity but concentrated on winning a new empire of style. Ramón María del Valle-Inclán (1866–1936) and Ramón Pérez de Ayala (1880–1962) brought a highly original lyricism to the novel, while Pío Baroja (1872–1956) concentrated on representing a world of discrete events, unbound by a unifying philosophy. The fiction that came out of the Civil War of 1936–39 returned to a kind of didactic realism, as with the novels of José María Gironella, who depicted a ravaged and suffering Spain. Camilo José Cela, perhaps the most important modern Spanish novelist, combines realism with a highly original style. His *Familia de Pascual Duarte*—harrowing, compassionate, brilliantly economical—is a novel of towering merit.

Italian. Though Italy originated the novella, it was slow in coming to the full-length novel. There is little to record before *I promessi sposi* (*The Betrothed*) by Alessandro Manzoni (1785–1873), a romantic and patriotic novel that describes life in Italy under Spanish domination in the 17th century. That combination of regionalism and realism already noted in the fiction of Germany and Spain did not appear in Italy until after the unification in 1870, when Giovanni Verga (1840–1922) celebrated his native Sicily in *I malavoglia* (*The House by the Medlar Tree*) and Antonio Fogazzaro (1842–1911) showed life in northern Italy during the struggle for unification in *Piccolo mondo antico* (*The Little World of the Past*). Gabriele D'Annunzio (1863–1938) and Luigi Pirandello (1867–1936) were too original for easy classification, and both worked in all the literary media. Pirandello, especially, helped to bring Italian literature into the modern world through such philosophical novels as *Il fu Mattia Pascal* (*The Late Matthew Pascal*), which raises profound questions about the nature of human identity and yet contrives to be witty, sunny, and essentially Italianate. The importance of Italo Svevo (1861–1928) was obscured for some time because of the difficulty of Italian literati in accepting his relatively unadorned style, but works like *La coscienza di Zeno* (1923; *Confessions of Zeno*) and *Senilità* (a title that James Joyce, Svevo's friend and English teacher, translated *As a Man Grows Older*) are generally recognized as major contributions to the international novel.

Pirandello

The significant fiction of the Mussolini regime was produced by anti-Fascist exiles like Giuseppe Borgese (1882–1952) and Ignazio Silone (1900–78), whose *Pane e vino* (*Bread and Wine*) is accepted as a 20th-century classic. Alberto Moravia, with *La romana* (*The Woman of Rome*) and *La noia* (*The Empty Canvas*), is perhaps the most popular Italian novelist outside Italy, but he is probably less important than Giuseppe Berto, Cesare Pavese (1908–50), and Elio Vittorini (1908–66). Giuseppe di Lampedusa (1896–1957) created a solitary masterpiece in *Il gattopardo* (*The Leopard*). The experimental writing of Italo Calvino and Carlo Emilio Gadda has become better known outside Italy, but the Italian novel remains linguistically conservative and needs the impact of some powerfully iconoclastic literary figure like James Joyce.

Scandinavian languages. Norway is better known for Ibsen's contribution to the drama and Grieg's to music than for fiction of the first quality. Nevertheless, the novels of Knut Hamsun (1859–1952) earned him the Nobel Prize in 1920, and Sigrid Undset (1882–1949) received the same honour in 1928, though the work of both has failed to engage the lasting attention of world readers of fiction. In Denmark, Johannes Vilhelm Jensen (1873–1950), another Nobel Prize winner, Isak Dinesen (the pen name of Baroness Karen Blixen, 1885–1962), and Martin

Nexö (1869–1954) have contributed to their country's fictional literature but made little mark beyond. Except for the Nobel Prize winners Eyvind Johnson (1900–76) and Harry Martinson (1904–78), the achievement of Sweden's novelists is inconsiderable. The novels of Halldór Laxness have restored to Iceland some of the literary fame it once earned for its sagas.

Slavic and East European languages. The greatest novelist of Czech origin, Franz Kafka, wrote in German, but writers in the vernacular include world-famous names such as Jarsoslav Hašek (1883–1923), whose *The Good Soldier Schweik* is acknowledged to be a comic masterpiece, and Karel Čapek (1890–1938), best known for the plays *R.U.R.* (which gave "robot" to the world's vocabulary) and *The Life of the Insects* but also notable for the novels *The Absolute at Large*, *Krakatit*, and *The War with the Newts*. The Czech fictional genius tends to the satirical and the fantastic. Serbian fiction gained international recognition in 1961 when Ivo Andrić won the Nobel Prize for his novels, notably those dealing with the history of Bosnia, which he had written during the second World War. Poland can claim two Nobel prizewinning novelists in Henryk Sienkiewicz (1846–1916), the author of *Quo Vadis?*, and Władysław S. Reymont (1868–1925), whose novel *The Peasants* is an aromatic piece of bucolic realism. Witold Gombrowicz (1904–69) wrote a novel, *Ferdydurke*, that was subjected to two separate modes of suppression—first Fascist, later Communist. A remarkable surrealist essay on "anal tyranny" and a depersonalization, it still awaits the acclaim that is due. Romania, which produced outstanding novelists in Eugen Lovinescu and Titu Maiorescu, suffered, like other Balkan countries, from the totalitarian suppression of the free creative spirit; but Dumitru Radu Popescu, author of *The Blue Lion*, was bold enough in the 1960s to question Communist orthodoxy through the medium of fiction. The work of the Hungarian Gusztáv Rab (1901–66) awaits recognition. His brilliant *Sabaria* develops very courageously the theme of the conflict between Communism and Christianity and reaches conclusions favourable to neither. It is a disturbing and beautifully composed book.

Czech
fiction

Modern Greece, perhaps more famous for its poets George Seferis and Constantine Cavafy, has produced at least one major novelist in Nikos Kazantzákis (1885–1957), whose *Zorba the Greek* became famous through its film adaptation. *The Last Temptation of Christ*, which presents the life of Jesus as a struggle to overcome "the dark immemorial forces of the Evil One, human and pre-human," glows with the writer's personality and bristles with the dialectic that was one of the first Greek gifts to Western civilization.

Kazantzákis

The Jewish novel. The literature of the Diaspora—the dispersion of the Jews after their exile from ancient Babylonia—records many large achievements in the languages of exile and that dialect of Low German—Yiddish—which the Ashkenazi Jews have taken around the world. Perhaps the most interesting of modern Jewish novelists in Yiddish is Isaac Bashevis Singer, a naturalized American who refuses to be absorbed linguistically into America, unlike Bellow and Malamud, who have brought to the Anglo-American language typical tones and rhythms of the ghetto.

Israel, which is producing its own rich crop of national writers, began with an existing corpus of literature in modern Hebrew, a language promoted and nurtured by such scholars as Eliezer ben Yehuda (1858–1923) and the members of the neologizing Hebrew Academy of Israel. Among the early Hebrew novelists is Abraham Mapu (1808–67); among the later ones is the brilliant Moshe Shamir. The first Nobel Prize for Literature ever awarded to a Hebrew novelist went, justly, to Shmuel Yosef Agnon (1888–1970). The contemporary Hebrew novel is notable for a fusion of international sophistication and earthy homegrown realism.

ASIA, AFRICA, LATIN AMERICA

China. The tradition of storytelling is an ancient one in China, and the full-length novel can be found as far back as the late 16th century. The fiction of the 18th century

shows a variety of themes and techniques not dissimilar to those of Europe, with social satire, chivalric romance, and adventure stories. Ts'ao Chan (1715?-63) wrote a novel called *Hung lou meng* (translated as *Dream of the Red Chamber* in 1892), which has features not unlike those of Galsworthy's *The Forsyte Saga* and Mann's *Buddenbrooks*—a story of a great aristocratic family in decline, garnished with love interest and shot through with pathos. The early days of the 20th century saw the foundation of the Chinese republic and the development of popular fiction written in the vernacular style of Chinese called *pai-hua*.

From 1917 until the Sino-Japanese War (1937-45), a period of social and intellectual ferment, there was a great influx of Western novels, and, under their influence, movements like those devoted to realism and naturalism in Europe produced a great number of didactic novels. Novels like Lao She's *Rickshaw Boy* and Pa Chin's *Chinese Earth* appeared in the postwar period, but the Communist regime has quelled all but propagandist fiction. Such novels as Liu Ching's *Wall of Bronze*, Chao Shu-li's *Changes in Li Village*, and Ting Ling's *Sun Shines over the Sangkan River* are typical glorifications of the Maoist philosophy and the socialist achievements of the people.

Japan. A literature that admires economy, like the Japanese, is bound to favour, in fictional art, the short story above the full-length novel. Nevertheless, some of the ancient pillow-books, with their diary jottings and anecdotes, have the ring of autobiographical novels; while Murasaki Shikibu's *Tale of Genji*, produced nearly a thousand years ago, is a great and sophisticated work of fictional history. The 20th-century Japanese novel has been developed chiefly under Western influence, like the work of Akutagawa Ryunosuke (1892-1927), whose short stories "Rashōmon" and "Yabu no naka" became the highly applauded film *Rashomon*. Tanizaki Jun-ichirō (1886-1965) is well known in America and Europe for *The Key*, *The Makioka Sisters*, and *Diary of a Mad Old Man*. His novels, all set in a modern, Westernized Japan, are assured of universal popularity because of their frank and lavish sexual content. Mishima Yukio (1925-70), who committed ceremonial suicide at the height of his reputation, was perhaps the most successful, and certainly the most prolific, of all modern Japanese novelists. Ten of his fictions had been filmed; he had won all the major Japanese literary awards and appeared to be destined for the Nobel Prize. His works are characterized by ruthless violence and a perversity that, however much it seems to derive from the pornographic excesses of Western fiction, is certainly in the Japanese tradition. His work is hard to judge, but he was the most considerable literary figure of the East.

India and East Asia. The Indian novel is a branch of British Commonwealth literature. That is to say, it is practiced by writers who have received a British literary education and, for the most part, publish their books in London. There is no evidence of any great development of fiction in any of the native Indian tongues: a taste for reading novels, as opposed to seeing films or reading short stories in the Hindi or Punjabi press, is acquired in India along with an education in English, which is still the unifying tongue of the subcontinent. In Malaysia, where the same tradition holds, short stories are being written in Malay, Chinese, and Tamil, but the full-length novel is almost exclusively written in English. It is as much a matter of literary markets as of literary education. No novel in any of the tongues of the peninsula is likely to be a remunerative publishing proposition, as Han Suyin—a best-selling Chinese woman doctor from Johore, who found large fame with *A Many-Splendoured Thing*—would be the first to admit, for all her dislike of the British colonial tradition. Indonesia's abandonment of all vestiges of its Dutch colonial past is associated with the encouragement of a literature in Bahasa Indonesia, a variety of Malay, and there are a number of Indonesian novelists still looking for a large educated audience—inevitably in translation, in the West. Among these Mochtar Lubis is notable; his *Twilight in Djakarta* is a bitter indictment of the Sukarno regime, which promptly sent him to prison. This was not the kind of fiction that the new Indonesia had in mind.

Africa. What applies to India and the East Indies applies also to Africa: there is still an insufficiently large audience for fiction written in any of the major African languages, and African novelists brought up on English are only too happy to continue working in it. Other European languages—chiefly Afrikaans (a South African variety of Dutch) and French—are employed as a fictional medium. Arthur Fula, who wrote *Janie Giet die Beeld* (*Janie Casts the Image*), is an Afrikaans novelist whose reputation has not, as yet, stretched to Europe or America, but the novelists of the former French possessions are gaining a name among serious students of the African novel. Mongo Beti, from Cameroon, is known for his *Pauvre Christ de Bomba* and *Le Roi miraculé*; the Ivory Coast has Aké Loba, and Hamidou Kane represents Senegal. This new African French deserves to be regarded as a distinct literary language, unrelated to that of Paris or Quebec, but the critical and linguistic tools for appraising the fiction of French-speaking Africans are not yet available.

Latin America. One of the greatest contemporary writers of fiction, Jorge Luis Borges, is an Argentinian, but his *ficciones* are very short short stories and must, with regrets, be excluded from any survey of the novel. It is significant, however, that the circumstances for the creation of a great fictional literature are in existence in Latin America, signifying (unlike the case in many former British dependencies) a shedding of the old colonial provincialism, which relied on the opinion of Madrid or, in Brazil, of Lisbon. The first Latin American novel was probably *El periquillo samiento* (1816; *The Itching Parrot*), by the Mexican José Joaquín Fernández de Lizardi (1776-1827), a picaresque work satirizing colonial conditions. In Argentina, José Mármol (1817-1871) published the first major novel of the continent—*Amalia* (1851-55), a powerful study of the fear and degradation that were rife in Buenos Aires during the dictatorship of the corrupt and tyrannical Juan Manuel de Rosas. The Romantic period in Europe had its counterpart in the sentimental wave that overtook such novelists as the Colombian Jorge Isaacs (1837-95), while the new humanitarianism found a voice in *Birds Without a Nest*, a protest-novel on the conditions forced on the Indians of Peru, written by Clorinda Matto de Turner (1854-1901). Juan León Mera followed the same trend in *Cumandá*, a novel about the oppression of the Ecuadorian Indians. As Latin America moved toward the modern age, the inevitable novels of urban social protest made their appearance. Alberto Blest Gana (1830-1920) of Chile, Carlos Reyles (1868-1938) of Uruguay, and Gustavo Martínez Zuviría of Argentina are names of some historical significance; and Reyles's naturalistic novel *La Raza de Caín* (*Cain's Race*) is original in that it finds a parallel between the breeding of stock and the building of a human society.

A 20th-century reaction against the bourgeois novel led to the movement known as nativism, with its concentration on the land itself, the lot of the indigenous peoples, the need for an anti-racist revolution with the aim of true egalitarianism. *Los de Abajo* (1915; *The Underdogs*), by the Mexican Mariano Azuela (1873-1952), and *El Señor Presidente* (1946), by Miguel Angel Asturias of Guatemala, typify the new revolutionary novel. Much of the didactic energy that went to the making of such work resulted in a lack of balance between subject matter and style: the literature of revolt tends to be shrill and crude. The inevitable reaction to a more sophisticated kind of literature, in which the individual became more important than society and the unconscious more interesting than the operation of reason, led to the highly refined experimentation that is the mark of Borges, the Brazilian Érico Veríssimo, and Eduardo Barrios of Chile.

The fiction of Brazil is in many ways more interesting than the fiction of the mother country, which has produced only one major novelist in the realist José Maria de Eça de Queirós (1845-1900), Gregório de Matos Guerra, as early as the 17th century, wrote bitterly of colonial administration and painted a realistic picture of Brazilian life. Irony and keen observation have been characteristic of the Brazilian novel ever since—as in the *Memórias Póstumas de Brás Cubas* (1881) and *Dom Casmurro* of

Western influences

Rise of nativism

Joaquim Maria Machado de Assis (1839–1908); the *Canaã* (1902; *Canaan*) of José Pereira da Graça Aranha (1868–1931), a remarkable study of the disillusion of German settlers in a new land; and the *Os Sertões* (*Rebellion in the Backlands*) of Euclides da Cunha (1866–1909), an account of a revolt against the newly formed Brazilian republic. The social consciousness of Brazilian novelists is remarkable, especially when it is associated with a strong concern for stylistic economy and grace. Monteiro Lobato (1882–1948), with his rustic hero Jeca Tatú; José Lins do Rêgo (1901–1957), chronicler of the decay of the old plantation life; Jorge Amado, a socialist novelist much concerned with slum life in Bahia; Érico Veríssimo, an experimental writer grounded in the traditional virtues of credible plot and strong characterization—these attest a vigour hardly to be found in the fiction published in Lisbon.

SOCIAL AND ECONOMIC ASPECTS

Though publishers of fiction recognize certain obligations to art, even when these are unprofitable (as they usually are), they are impelled for the most part to regard the novel as a commercial property and to be better pleased with large sales of indifferent work than with the mere unremunerative acclaim of the intelligentsia for books of rare merit. For this reason, any novelist who seeks to practice his craft professionally must consult the claims of the market and effect a compromise between what he wishes to write and what the public will buy. Many worthy experimental novels, or novels more earnest than entertaining, gather dust in manuscript or are circulated privately in photocopies. Indeed, the difficulty that some unestablished novelists find in gaining a readership (which means the attention of a commercial publisher) has led them to take the copying machine as seriously as the printing press and to make the composition, mimeographing, binding, and distribution of a novel into a single cottage industry. For the majority of novelists the financial rewards of their art are nugatory, and only a strong devotion to the form for its own sake can drive them to the building of an oeuvre. The subsidies provided by university sinecures sustain a fair number of major American novelists; others, in most countries, support their art by practicing various kinds of subliterate—journalism, film scripts, textbooks, even pseudonymous pornography. Few novelists write novels and novels only.

Awards, patronage, and ancillary activities

There are certain marginal windfalls, and the hope of gaining one of these tempers the average novelist's chronic desperation. America has its National Book Award as well as its book club choices; France has a great variety of prizes; there are also international bestowals; above all there glows the rarest and richest of all accolades—the Nobel Prize for Literature. Quite often the Nobel Prize winner needs the money as much as the fame, and his election to the honour is not necessarily a reflection of a universal esteem which, even for geniuses like Samuel Beckett, means large sales and rich royalties. When Sinclair Lewis received the award in 1930, wealth and fame were added to wealth and fame already sufficiently large; when William Faulkner was chosen in 1949, most of his novels had been long out of print in America.

Prizes come so rarely, and often seem to be bestowed so capriciously, that few novelists build major hopes on them. They build even fewer hopes on patronage: Harriet Shaw Weaver, James Joyce's patroness, was probably the last of a breed that, from Maecenas on, once intermittently flourished; state patronage—as represented, for instance, by the annual awards of the Arts Council of Great Britain—can provide little more than a temporary palliative for the novelist's indigence. Novelists have more reasonable hopes from the world of the film or the stage, where adaptations can be profitable and even salvatory. The long struggles of the British novelist T.H. White came to an end when his Arthurian sequence *The Once and Future King* (1958) was translated into a stage musical called *Camelot*, though, by treating the lump sum paid to him as a single year's income instead of a reward for decades of struggle, nearly all the windfall would have gone for taxes if White had not taken his money into low-tax exile. Such writers as Graham Greene, nearly all

of whose novels have been filmed, must be tempted to regard mere book sales as an inconsiderable aspect of the rewards of creative writing. There are few novelists who have not received welcome and unexpected advances on film options, and sometimes the hope of film adaptation has influenced the novelist's style. In certain countries, such as Great Britain but not the United States, television adaptation of published fiction is common, though it pays the author less well than commercial cinema.

When a novelist becomes involved in film-script writing—either the adaptation of his own work or that of others—the tendency is for him to become subtly corrupted by what seems to him an easier as well as more lucrative technique than that of the novel. Most novelists write dialogue with ease, and their contribution to a film is mostly dialogue: the real problem in novel writing lies in the management of the *récit*. A number of potentially fine novelists, like Terry Southern and Frederick Raphael, have virtually abandoned the literary craft because of their continued success with script writing. In 70-odd years the British novelist Richard Hughes produced only three novels, the excellence of which has been universally recognized; fiction lovers have been deprived of more because of the claims of the film world on Hughes's talent. This kind of situation finds no counterpart in any other period of literary history, except perhaps in the Elizabethan, when the commercial lure of the drama made some good poets write poor plays.

The majority of professional novelists must look primarily to book sales for their income, and they must look decreasingly to hardcover sales. The novel in its traditional format, firmly stitched and sturdily clothbound, is bought either by libraries or by readers who take fiction seriously enough to wish to acquire a novel as soon as it appears: if they wait 12 months or so, they can buy the novel in paper covers for less than its original price. This edition of a novel has become, for the vast majority of fiction readers, the form in which they first meet it, and the novelist who does not achieve paperback publication is missing a vast potential audience. He may not repine at this, since the quantitative approach to literary communication may safely be disregarded: the legend on a paperback cover—FIVE MILLION COPIES SOLD—says nothing about the worth of the book within. Nevertheless, the advance he will receive from his hardcover publisher is geared to eventual paperback expectations, and the "package deal" has become the rule in negotiations between publisher and author's agent. The agent, incidentally, has become important to both publisher and author to an extent that writers like Daniel Defoe and Samuel Richardson would, if resurrected, find hard to understand.

The novelist may reasonably expect to augment his income through the sale of foreign rights in his work, though the rewards accruing from translation are always uncertain. The translator himself is usually a professional and demands a reasonable reward for his labours, more indeed than the original author may expect: the reputations of some translators are higher than those of some authors, and even the translators' names may be better known. Moreover, the author who earns most from publication in his own language will usually earn most in translation, since it is the high initial home sales that attract foreign publishers to a book. The more "literary" a novel is, the more it exploits the resources of the author's own language, the less likely is it to achieve either popularity at home or publication abroad. Best-selling novels like Mario Puzo's *Godfather* (1969) or Arthur Hailey's *Airport* (1968) are easy to read and easy to translate, so they win all around. It occasionally happens that an author is more popular abroad than he is at home: the best-selling novels of the Scottish physician-novelist A.J. Cronin are no longer highly regarded in England and America, as they were in the 1930s and '40s, but they continued to sell by the million in the U.S.S.R. several decades later. However, a novelist is wisest to expect most from his own country and to regard foreign popularity as an inexplicable bonus.

As though his financial problems were not enough, the novelist frequently has to encounter those dragons unleashed by public morality or by the law. The struggles

The importance of the paperback

Problems of obscenity and libel

of Flaubert, Zola, and Joyce, denounced for attempting to advance the frontiers of literary candour, are well known and still vicariously painful, but lesser novelists, working in a more permissive age, can record cognate agonies. Generally speaking, any novelist writing after the publication in the 1960s of Hubert Selby's *Last Exit to Brooklyn* or Gore Vidal's *Myra Breckenridge* can expect little objection, on the part of either publisher or police, to language or subject matter totally unacceptable, under the obscenity laws then operating, in 1922, when *Ulysses* was first published. This is certainly true of America, if not of Ireland or Malta. But many serious novelists fear an eventual reaction against literary permissiveness as a result of the exploitation by cynical obscenity mongers or hard-core pornographers of the existing liberal situation.

In some countries, particularly Great Britain, the law of libel presents insuperable problems to novelists who, innocent of libellous intent, are nevertheless sometimes charged with defamation by persons who claim to be the models for characters in works of fiction. Disclaimers to the effect that "resemblances to real-life people are wholly coincidental" have no validity in law, which upholds the right of a plaintiff to base his charge on the corroboration of "reasonable people." Many such libel cases are settled before they come to trial, and publishers will, for the sake of peace and in the interests of economy, make a cash payment to the plaintiff without considering the author's side. They will also, and herein lies the serious blow to the author, withdraw copies of the allegedly offensive book and pulp the balance of a whole edition. Novelists are seriously hampered in their endeavours to show, in a traditional spirit of artistic honesty, corruption in public life; they have to tread carefully even in depicting purely imaginary characters and situations, since the chance collocation of a name, a profession, and a locality may produce a libellous situation.

EVALUATION AND STUDY

It has been only in comparatively recent times that the novel has been taken sufficiently seriously by critics for the generation of aesthetic appraisal and the formulation of fictional theories. The first critics of the novel developed their craft not in full-length books but in reviews published in periodicals: much of this writing—in the late 18th and early 19th centuries—was of an occasional nature, and not a little of it casual and desultory; nor, at first, did critics of fiction find it easy to separate a kind of moral judgment of the subject matter from an aesthetic judgement of the style. Such fragmentary observations on the novel as those made by Dr. Johnson in conversation or by Jane Austen in her letters, or, in France, by Gustave Flaubert during the actual process of artistic gestation, have the charm and freshness of insight rather than the weight of true aesthetic judgment. It is perhaps not until the beginning of the 20th century, when Henry James wrote his authoritative prefaces to his own collected novels, that a true criteriology of fiction can be said to have come into existence. The academic study of the novel presupposes some general body of theory, like that provided by Percy Lubbock's *Craft of Fiction* (1921) or E.M. Forster's *Aspects of the Novel* (1927) or the subsequent writings of the critics Edmund Wilson and F.R. Leavis. Since World War II it may be said that university courses in the evaluation of fiction have attained the dignity traditionally monopolized by poetry and the drama.

A clear line should be drawn between the craft of fiction criticism and the journeyman work of fiction reviewing. Reviews are mainly intended to provide immediate information about new novels: they are done quickly and are subject to the limitations of space; they not infrequently make hasty judgments that are later regretted. The qualifications sought in a reviewer are not formidable: smartness, panache, waspishness—qualities that often draw the attention of the reader to the personality of the reviewer rather than the work under review—will always be more attractive to circulation-hunting editors than a less spectacular concern with balanced judgment. A thoughtful editor will sometimes put the reviewing of novels into the hands of a practicing novelist, who—knowing the labour that goes

into even the meanest book—will be inclined to sympathy more than to flamboyant condemnation. The best critics of fiction are probably novelists *manqués*, men who have attempted the art and, if not exactly failed, not succeeded as well as they could have wished. Novelists who achieve very large success are possibly not to be trusted as critics: obsessed by their own individual aims and attainments, shorn of self-doubt by the literary world's acclaim or their royalty statements, they bring to other men's novels a kind of magisterial blindness.

Novelists can be elated by good reviews and depressed by bad ones, but it is rare that a novelist's practice is much affected by what he reads about himself in the literary columns. Genuine criticism is a very different matter, and a writer's approach to his art can be radically modified by the arguments and summations of a critic he respects or fears. As the hen is unable to judge of the quality of the egg it lays, so the novelist is rarely able to explain or evaluate his work. He relies on the professional critic for the elucidation of the patterns in his novels, for an account of their subliminal symbolism, for a reasoned exposition of their stylistic faults. As for the novel reader, he will often learn enthusiasm for particular novelists through the writings of critics rather than from direct confrontation with the novels themselves. The essays in Edmund Wilson's *Axel's Castle* (1931) aroused an interest in the Symbolist movement which the movement was not easily able to arouse by itself; the essay on *Finnegans Wake*, collected in Wilson's *Wound and the Bow* (1941), eased the way into a very difficult book in a manner that no grim work of solid exegesis could have achieved. The essence of the finest criticism derives from wisdom and humanity more than from mere expert knowledge. Great literature and great criticism possess in common a sort of penumbra of wide but unsystematic learning, a devotion to civilized values, an awareness of tradition, and a willingness to rely occasionally on the irrational and intuitive.

All this probably means that the criticism of fiction can never, despite the efforts of aestheticians schooled in modern linguistics, become an exact science. A novel must be evaluated in terms of a firmly held literary philosophy, but such a philosophy is, in the final analysis, based on the irrational and subjective. If the major premises on which F.R. Leavis bases his judgments of George Eliot, Mark Twain, and D.H. Lawrence are accepted, then an acceptance of the judgments themselves is inescapable. But many students of fiction who are skeptical of Leavis will read him in order that judgments of their own may emerge out of a purely negative rejection of his. In reading criticism a kind of dialectic is involved, but no synthesis is ever final. The process of reevaluation goes on for ever. One of the sure tests of a novel's worth is its capacity for engendering critical dialectic: no novel is beyond criticism, but many are beneath it.

THE FUTURE OF THE NOVEL

It is apparent that neither law nor public morality nor the public's neglect nor the critic's scorn has ever seriously deflected the dedicated novelist from his self-imposed task of interpreting the real world or inventing alternative worlds. Statistics since World War II have shown a steady increase in the number of novels published annually, and beneath the iceberg tip of published fiction lies a submarine Everest of unpublished work. It has been said that every person has at least one novel in him, and the near-universal literacy of the West has produced dreams of authorship in social ranks traditionally deprived of literature. Some of these dreams come true, and taxi drivers, pugilists, criminals, and film stars have competed, often successfully, in a field that once belonged to professional writers alone. It is significant that the amateur who dreams of literary success almost invariably chooses the novel, not the poem, essay, or autobiography. Fiction requires no special training and can be readable, even absorbing, when it breaks the most elementary rules of style. It tolerates a literary incompetence unthinkable in the poem. If all professional novelists withdrew, the form would not languish: amateurs would fill the market with first and only novels, all of which would find readership.

The distinction between reviewing and criticism

Amateur novelists

But the future of any art lies with its professionals. Here a distinction has to be made between the Joyces, Henry Jameses, and Conrads on the one hand, and the more ephemeral Mickey Spillanes, Harold Robbinses, and Irving Wallaces on the other. Of the skill of the latter class of novelists there can be no doubt, but it is a skill employed for limited ends, chiefly the making of money, and through it the novel can never advance as art. The literary professionals, however, are dedicated to the discovery of new means of expressing, through the experiential immediacies that are the very stuff of fiction, the nature of man and society. In the symbiosis of publishing, the best-seller will probably continue to finance genuine fictional art. Despite the competition from other art media, and the agonies and the indigence, there are indications that the serious novel will flourish in the future.

It will flourish because it is the one literary form capable of absorbing all the others. The technique of the stage drama or the film can be employed in the novel (as in *Ulysses* and *Giles Goat-Boy*), as can the devices of poetry (as in Philip Toynbee's *Pantaloon* and the novels of Wilson Harris and Janet Frame). In France, as Michel Butor has pointed out, the new novel is increasingly performing some of the tasks of the old essay; in America, as Capote's *In Cold Blood* and Mailer's *Armies of the Night* have shown, the documentary report can gain strength from its presentation as fictional narrative. There are few limits on what the novel can do, there are many experimental paths still to be trod, and there is never any shortage of subject matter.

For all this, periods of decline and inanition may be expected, though not everywhere at once. The strength of the American novel in the period after World War II had something to do with the national atmosphere of breakdown and change: political and social urgencies promoted a quality of urgency in the works of such writers as Mailer, Bellow, Ellison, Heller, and Philip Roth. In the same period, Britain, having shed its empire and erected a welfare state, robbed its novelists of anything larger to write about than temporary indentations in the class system, suburban adultery, and manners. An achieved or static society does not easily produce great art. France, which has known much social and ideological turmoil, has generated a new aesthetic of the novel as well as a philosophy that, as Sartre and Camus have shown, is very suitable for fictional expression. A state on which intellectual quietism or a political philosophy of art is imposed by the ruling party can, as the Soviet Union and China show, succeed only in thwarting literary greatness, but the examples of Pasternak and Solzhenitsyn are reminders that repression can, with rare artistic spirits, act as an agonizing stimulus.

Every art in every country is subject to a cyclical process: during a period of decline it is necessary to keep the communication lines open, producing minor art so that it may some day, unexpectedly, turn into major art. Wherever the novel seems to be dying it is probably settling into sleep; elsewhere it will be alive and vigorous enough. It is important to believe that the novel has a future, though not everywhere at once. (An.B.)

Short story

The short story is a kind of prose fiction, usually more compact and intense than the novel and the short novel (novelette). Prior to the 19th century it was not generally regarded as a distinct literary form. But although in this sense it may seem to be a uniquely modern genre, the fact is that short prose fiction is nearly as old as language itself. Throughout history man has enjoyed various types of brief narratives: jests, anecdotes, studied digressions, short allegorical romances, moralizing fairy tales, short myths, and abbreviated historical legends. None of these constitutes a short story as the 19th and 20th centuries have defined the term, but they do make up a large part of the milieu from which the modern short story emerged.

Many of the elements of storytelling common to the short story and the novel are discussed at greater length in the preceding section on the novel. The short stories of particular literary cultures, along with other genres, are

discussed in articles such as LITERATURE, THE HISTORY OF WESTERN; and in articles on the arts of various peoples—e.g., SOUTH ASIAN ARTS.

ANALYSIS OF THE GENRE

As a genre, the short story has received relatively little critical attention, and the most valuable studies of the form that exist are often limited by region or era (e.g., Ray B. West's *The Short Story in America, 1900-50*). One recent attempt to account for the genre has been offered by the Irish short story writer Frank O'Connor, who suggests that stories are a means for "submerged population groups" to address a dominating community. Most other theoretical discussions, however, are predicated in one way or another on Edgar Allan Poe's thesis that stories must have a compact, unified effect.

By far the majority of criticism on the short story focuses on techniques of writing. Many, and often the best of the technical works, advise the young reader—alerting him to the variety of devices and tactics employed by the skilled writer. On the other hand, many of these works are no more than treatises on "how to write stories" for the young writer, and not serious critical material.

The prevalence in the 19th century of two words, "sketch" and "tale," affords one way of looking at the genre. In the United States alone there were virtually hundreds of books claiming to be collections of sketches (Washington Irving's *Sketch Book*, William Dean Howells' *Suburban Sketches*) or collections of tales (Poe's *Tales of the Grotesque and Arabesque*, Herman Melville's *Piazza Tales*). These two terms establish the polarities of the milieu out of which the modern short story grew.

The tale is much older than the sketch. Basically, the tale is a manifestation of a culture's unaging desire to name and conceptualize its place in the cosmos. It provides a culture's narrative framework for such things as its vision of itself and its homeland or for expressing its conception of its ancestors and its gods. Usually filled with cryptic and uniquely deployed motifs, personages, and symbols, tales are frequently fully understood only by members of the particular culture to which they belong. Simply, tales are intracultural. Seldom created to address an outside culture, a tale is a medium through which a culture speaks to itself and thus perpetuates its own values and stabilizes its own identity. The old speak to the young through tales.

The sketch, by contrast, is intercultural, depicting some phenomenon of one culture for the benefit or pleasure of a second culture. Factual and journalistic, in essence the sketch is generally more analytic or descriptive and less narrative or dramatic than the tale. Moreover, the sketch by nature is *suggestive*, incomplete; the tale is often *hyperbolic*, overstated.

The primary mode of the sketch is written: that of the tale, spoken. This difference alone accounts for their strikingly different effects. The sketch writer can have, or pretend to have, his eye on his subject. The tale, recounted at court or campfire—or at some place similarly removed in time from the event—is nearly always a recreation of the past. The tale-teller is an agent of *time*, bringing together a culture's past and its present. The sketch writer is more an agent of *space*, bringing an aspect of one culture to the attention of a second.

It is only a slight oversimplification to suggest that the tale was the only kind of short fiction until the 16th century, when a rising middle class interest in social realism on the one hand and in exotic lands on the other put a premium on sketches of subcultures and foreign regions. In the 19th century certain writers—those one might call the "fathers" of the modern story: Nikolay Gogol, Hawthorne, E.T.A. Hoffmann, Heinrich von Kleist, Prosper Mérimée, Poe—combined elements of the tale with elements of the sketch. Each writer worked in his own way, but the general effect was to mitigate some of the fantasy and stultifying conventionality of the tale and, at the same time, to liberate the sketch from its bondage to strict factuality. The modern short story, then, ranges between the highly imaginative tale and the photographic sketch and in some ways draws on both.

The short stories of Ernest Hemingway, for example,

The sketch and the tale

Modern fusions

may often gain their force from an exploitation of traditional mythic symbols (water, fish, groin wounds), but they are more closely related to the sketch than to the tale. Indeed, Hemingway was able at times to submit his apparently factual stories as newspaper copy. In contrast, the stories of Hemingway's contemporary William Faulkner more closely resemble the tale. Faulkner seldom seems to understate, and his stories carry a heavy flavour of the past. Both his language and his subject matter are rich in traditional material. A Southerner might well suspect that only a reader steeped in sympathetic knowledge of the traditional South could fully understand Faulkner. Faulkner may seem, at times, to be a Southerner speaking to and for Southerners. But, as, by virtue of their imaginative and symbolic qualities, Hemingway's narratives are more than journalistic sketches, so, by virtue of their explorative and analytic qualities, Faulkner's narratives are more than Southern tales.

Whether or not one sees the modern short story as a fusion of sketch and tale, it is hardly disputable that today the short story is a distinct and autonomous, though still developing, genre.

HISTORY

Origins. The evolution of the short story first began before man could write. To aid himself in constructing and memorizing tales, the early storyteller often relied on stock phrases, fixed rhythms, and rhyme. Consequently, many of the oldest narratives in the world, such as the famous Babylonian tale the *Epic of Gilgamesh* (c. 2000 BC), are in verse. Indeed, most major stories from the ancient Middle East were in verse: "The War of the Gods," "The Story of Adapa" (both Babylonian), "The Heavenly Bow," and "The King Who Forgot" (both Canaanite). These tales were inscribed in cuneiform on clay during the 2nd millennium BC.

The earliest tales extant from Egypt were composed on papyrus at a comparable date. The ancient Egyptians seem to have written their narratives largely in prose, apparently reserving verse for their religious hymns and working songs. One of the earliest surviving Egyptian tales, "The Shipwrecked Sailor" (c. 2000 BC), is clearly intended to be a consoling and inspiring story to reassure its aristocratic audience that apparent misfortune can in the end become good fortune. Also recorded during the 12th dynasty were the success story of the exile Sinuhe and the moralizing tale called "King Cheops [Khufu] and the Magicians." The provocative and profusely detailed story "The Tale of Two Brothers" (or "Anpu and Bata") was written down during the New Kingdom, probably around 1250 BC. Of all the early Egyptian tales, most of which are baldly didactic, this story is perhaps the richest in folk motifs and the most intricate in plot.

The earliest tales from India are not as old as those from Egypt and the Middle East. The *Brāhmaṇas* (c. 700 BC) function mostly as theological appendixes to the Four Vedas, but a few are composed as short, instructional parables. Perhaps more interesting as stories are the later tales in the Pāli language, *The Jātaka*. Although these tales have a religious frame that attempts to recast them as Buddhist ethical teachings, their actual concern is generally with secular behaviour and practical wisdom. Another, nearly contemporaneous collection of Indian tales, *The Pañca-tantra* (c. AD 500), has been one of the world's most popular books. This anthology of amusing and moralistic animal tales, akin to those of "Aesop" in Greece, was translated into Middle Persian in the 6th century; into Arabic in the 8th century; and into Hebrew, Greek, and Latin soon thereafter. Sir Thomas North's English translation appeared in 1570. Another noteworthy collection is *Kathā-saritsāgara* ("Ocean of Rivers of Stories), a series of tales assembled and recounted in narrative verse in the 11th century by the Sanskrit writer Samadeva. Most of these tales come from much older material, and they vary from the fantastic story of a transformed swan to a more probable tale of a loyal but misunderstood servant.

During the 2nd, 3rd, and 4th centuries BC, the Hebrews first wrote down some of their rather sophisticated narratives, which are now a part of the Old Testament and the

Apocrypha. The book of Tobit displays an unprecedented sense of ironic humour: Judith creates an unrelenting and suspenseful tension as it builds to its bloody climax; the story of Susanna, the most compact and least fantastic in the Apocrypha, develops a three-sided conflict involving the innocent beauty of Susanna, the lechery of the elders, and the triumphant wisdom of Daniel. The Old Testament books of Ruth, Esther, and Jonah hardly need mentioning: they may well be the most famous stories in the world.

Nearly all of the ancient tales, whether from Israel, India, Egypt, or the Middle East, were fundamentally didactic. Some of these ancient stories preached by presenting an ideal for readers to imitate. Others tagged with a "moral" were more direct. Most stories, however, preached by illustrating the success and joy that was available to the "good" man and by conveying a sense of the terror and misery that was in store for the wayward.

The early Greeks contributed greatly to the scope and art of short fiction. As in India, the moralizing animal fable was a common form; many of these tales were collected as "Aesop's fables" in the 6th century BC. Brief mythological stories of the gods' adventures in love and war were also popular in the pre-Attic age. Apollodorus of Athens compiled a handbook of epitomes, or abstracts, of these tales around the 2nd century BC, but the tales themselves are no longer extant in their original form. They appear, though somewhat transformed, in the longer poetical works of Hesiod, Homer, and the tragedians. Short tales found their way into long prose forms as well, as in Hellanicus' *Persika* (5th century BC, extant only in fragments).

Herodotus, the "father of history," saw himself as a maker and reciter of *logoi* (things for telling, tales). His long *History* is interspersed with such fictionalized digressions as the stories of Polycrates and his emerald ring, of Candaules' attractive wife, and of Rhampsinitus' stolen treasure. Xenophon's philosophical history, the *Cyropaedia* (4th century BC), contains the famous story of the soldier Abradates and his lovely and loyal wife Panthea, perhaps the first Western love story. The *Cyropaedia* also contains other narrative interpolations: the story of Phraules, who freely gave away his wealth; the tale of Gobryas' murdered son; and various anecdotes describing the life of the Persian soldier.

Moreover, the Greeks are usually credited with originating the romance, a long form of prose fiction with stylized plots of love, catastrophe, and reunion. The early Greek romances frequently took shape as a series of short tales. The *Love Romances* of Parthenius of Nicaea, who wrote during the reign of Augustus Caesar, is a collection of 36 prose stories of unhappy lovers. *The Milesian Tales* (no longer extant) was an extremely popular collection of erotic and ribald stories composed by Aristides of Miletus in the 2nd century BC and translated almost immediately into Latin. As the variety of these short narratives suggests, the Greeks were less insistent than earlier cultures that short fiction be predominantly didactic.

By comparison the contribution of the Romans to short narrative was small. Ovid's long poem, *Metamorphoses*, is basically a reshaping of over 100 short, popular tales into a thematic pattern. The other major fictional narratives to come out of Rome are novel-length works by Petronius (*Satyricon*, 1st century AD) and Apuleius (*The Golden Ass*, 2nd century AD). Like Ovid these men used potential short-story material as episodes within a larger whole. The Roman love of rhetoric, it seems, encouraged the development of longer and more comprehensive forms of expression. Regardless, the trend away from didacticism inaugurated by the Greeks was not reversed.

Middle Ages, Renaissance, and after. *Proliferation of forms.* The Middle Ages was a time of the proliferation, though not necessarily the refinement, of short narratives. The short tale became an important means of diversion and amusement. From the Dark Ages to the Renaissance, various cultures adopted short fiction for their own purposes. Even the aggressive, grim spirit of the invading Germanic barbarians was amenable to expression in short prose. The myths and sagas extant in Scandinavia and Iceland indicate the kinds of bleak and violent tales the invaders took with them into southern Europe.

The didacticism of early tales

The earliest tales extant

Beginnings of the romance

In contrast, the romantic imagination and high spirits of the Celts remained manifest in their tales. Wherever they appeared—in Ireland, Wales, or Brittany—stories steeped in magic and splendour also appeared. This spirit, easily recognized in such Irish mythological tales as *Longes mac n-Uislem* (probably 9th-century), infused the chivalric romances that developed somewhat later on the Continent. The romances usually addressed one of three “Matters”: the “Matter of Britain” (stories of King Arthur and his knights), the “Matter of France” (the Charlemagne cycle), or the “Matter of Rome” (stories out of antiquity, such as “Pyramus and Thisbe,” “Paris and Helen”). Many, but not all, of the romances are too long to be considered short stories. Two of the most influential contributors of short material to the “Matter of Britain” in the 12th century were Chrétien de Troyes and Marie de France. The latter was gifted as a creator of the short narrative poems known as the Breton lays. Only occasionally did a popular short romance like *Aucassin and Nicolette* (13th century) fail to address any of the three Matters.

Exempla
and
fabliaux

Also widely respected was the exemplum, a short didactic tale usually intended to dramatize or otherwise inspire model behaviour. Of all the exempla, the best known in the 11th and 12th centuries were the lives of the saints, some 200 of which are extant. The *Gesta Romanorum* (“Deeds of the Romans”) offered skeletal plots of exempla that preachers could expand into moralistic stories for use in their sermons.

Among the common people of the late Middle Ages there appeared a literary movement counter to that of the romance and exemplum. Displaying a preference for common sense, secular humour, and sensuality, this movement accounted in a large way for the practical-minded animals in beast fables, the coarse and “merry” jestbooks, and the ribald fabliaux. All were important as short narratives, but perhaps the most intriguing of the three are the fabliaux. First appearing around the middle of the 12th century, fabliaux remained popular for 200 years, attracting the attention of Boccaccio and Chaucer. Some 160 fabliaux are extant, all in verse.

The
framing
circum-
stance

Often, the medieval storyteller—regardless of the kind of tale he preferred—relied on a framing circumstance that made possible the juxtaposition of several stories, each of them relatively autonomous. Since there was little emphasis on organic unity, most storytellers preferred a flexible format, one that allowed tales to be added or removed at random with little change in effect. Such a format is found in *The Seven Sages of Rome*, a collection of stories so popular that nearly every European country had its own translation. The framing circumstance in *The Seven Sages* involves a prince condemned to death; his advocates (the seven sages) relate a new story each day, thereby delaying the execution until his innocence is made known. This technique is clearly similar to that of *The Arabian Nights*, another collection to come out of the Middle Ages. The majority of the stories in *The Arabian Nights* are framed by the story of Scheherazade in “A Thousand and One Nights.” Records indicate that the basis of this framing story was a medieval Persian collection, *Hezar Efsan* (“Thousand Romances,” no longer extant). In both the Persian and Arabian versions of the frame, the clever Scheherazade avoids death by telling her king-husband a thousand stories. Though the framing device is identical in both versions, the original Persian stories within the frame were replaced or drastically altered as the collection was adapted by the Arabs during the Muslim Manlūk period (AD 1250–1517).

Boccaccio
and
Chaucer

Refinement. Short narrative received its most refined treatment in the Middle Ages from Chaucer and Boccaccio. Chaucer’s versatility reflects the versatility of the age. In “The Miller’s Tale” he artistically combines two fabliaux; in “The Nun’s Priest’s Tale” he draws upon material common to beast fables; in “The Pardoner’s Tale” he creates a brilliantly revealing sermon, complete with a narrative exemplum. This short list hardly exhausts the catalogue of forms Chaucer experimented with. By relating tale to teller and by exploiting relationships among the various tellers, Chaucer endowed *The Canterbury Tales* with a unique, dramatic vitality.

Boccaccio’s genius, geared more toward narrative than drama, is of a different sort. Where Chaucer reveals a character through actions and assertions, Boccaccio seems more interested in stories as pieces of action. With Boccaccio, the characters telling the stories, and usually the characters within, are of subordinate interest. Like Chaucer, Boccaccio frames his well-wrought tales in a metaphoric context. The trip to the shrine at Canterbury provides a meaningful backdrop against which Chaucer juxtaposes his earthy and pious characters. The frame of the *Decameron* (from the Greek *deka*, 10, and *hēmera*, day) has relevance as well: during the height of the Black Plague in Florence, Italy, 10 people meet and agree to amuse and divert each other by telling 10 stories each. Behind every story, in effect, is the inescapable presence of the Black Death. The *Decameron* is fashioned out of a variety of sources, including fabliaux, exempla, and short romances.

Spreading popularity. Immediately popular, the *Decameron* produced imitations nearly everywhere. In Italy alone, there appeared at least 50 writers of *novelle* (as short narratives were called) after Boccaccio.

Learning from the success and artistry of Boccaccio and, to a lesser degree, his contemporary Franco Sacchetti, Italian writers for three centuries kept the Western world supplied with short narratives. Sacchetti was no mere imitator of Boccaccio. More of a frank and unadorned realist, he wrote—or planned to write—300 stories (200 of the *Trecentonovelle* [“300 Short Stories”] are extant) dealing in a rather anecdotal way with ordinary Florentine life. Two other well-known narrative writers of the 14th century, Giovanni Fiorentino and Giovanni Sercambi, freely acknowledged their imitation of Boccaccio. In the 15th century Masuccio Salernitano’s collection of 50 stories, *Il novellino* (1475), attracted much attention. Though verbosity often substitutes for eloquence in Masuccio’s stories, they are witty and lively tales of lovers and clerics.

Italian
short
narratives

With Masuccio the popularity of short stories was just beginning to spread. Almost every Italian in the 16th century, it has been suggested, tried his hand at *novelle*. Matteo Bandello, the most influential and prolific writer, attempted nearly everything from brief histories and anecdotes to short romances, but he was most interested in tales of deception. Various other kinds of stories appeared. Agnolo Firenzuolo’s popular *Ragionamenti d’amore* (“The Reasoning of Love”) is characterized by a graceful style unique in tales of ribaldry; Anton Francesco Doni included several tales of surprise and irony in his miscellany, *I marmi* (“The Marbles”); and Gianfrancesco Straparola experimented with common folktales and with dialects in his collection, *Le piacevoli notti* (“The Pleasant Nights”). In the early 17th century, Giambattista Basile attempted to infuse stock situations (often of the fairy-tale type, such as “Puss and Boots”) with realistic details. The result was often remarkable—a tale of hags or princes with very real motives and feelings. Perhaps it is the amusing and diverting nature of Basile’s collection of 50 stories that has reminded readers of Boccaccio. Or, it may be his use of a frame similar to that in the *Decameron*. Whatever the reason, Basile’s *Cunto de li cunti* (1634; *The Story of Stories*) is traditionally linked with Boccaccio and referred to as *The Pentamerone* (“The Five Days”). Basile’s similarities to Boccaccio suggest that in the 300 years between them the short story may have gained repute and circulation, but its basic shape and effect hardly changed.

This pattern was repeated in France, though the impetus provided by Boccaccio was not felt until the 15th century. A collection of 100 racy anecdotes, *Les Cent Nouvelles Nouvelles*, “The Hundred New Short Stories” (c. 1460), outwardly resembles the *Decameron*. Margaret of Angoulême’s *Heptaméron* (1558–59; “The Seven Days”), an unfinished collection of 72 amorous tales, admits a similar indebtedness.

In the early 17th century Béroalde de Verville placed his own Rabelaisian tales within a banquet frame in a collection called *Le Moyen de parvenir*, “The Way of Succeeding” (c. 1610). Showing great narrative skill, Béroalde’s stories are still very much in the tradition of Boccaccio; as a collection of framed stories, their main intent is to amuse and divert the reader.

As the most influential nation in Europe in the 15th and 16th centuries, Spain contributed to the proliferation of short prose fiction. Especially noteworthy are: Don Juan Manuel's collection of lively exempla *Libro de los enxiemplos del conde Lucanor et de Patronio* (1328–35), which antedates the *Decameron*; the anonymous story "The Abencerraje," which was interpolated into a pastoral novel of 1559; and, most importantly, Miguel de Cervantes' experimental *Novelas ejemplares* (1613: "Exemplary Novels"). Cervantes' short fictions vary in style and seriousness, but their single concern is clear: to explore the nature of man's secular existence. This focus was somewhat new for short fiction, heretofore either didactic or escapist.

Despite the presence of these and other popular collections, short narrative in Spain was eventually overshadowed by a new form that began to emerge in the 16th century—the novel. Like the earlier Romans, the Spanish writers of the early Renaissance often incorporated short story material as episodes in a larger whole.

Decline of short fiction. The 17th and 18th centuries mark the temporary decline of short fiction. The causes of this phenomenon are many: the emergence of the novel; the failure of the Boccaccio tradition to produce in three centuries much more than variations or imitations of older, well-worn material; and a renaissance fascination with drama and poetry, the superior forms of classical antiquity. Another cause for the disappearance of major works of short fiction is suggested by the growing preference for journalistic sketches. The increasing awareness of other lands and the growing interest in social conditions (accommodated by a publication boom) produced a plethora of descriptive and biographical sketches. Although these journalistic elements later were incorporated in the fictional short story, for the time being fact held sway over the imagination. Travel books, criminal biographies, social description, sermons, and essays occupied the market. Only occasionally did a serious story find its way into print, and then it was usually a production of an established writer like Voltaire or Addison.

Perhaps the decline is clearest in England, where the short story had its least secure foothold. It took little to obscure the faint tradition established in the 16th and 17th centuries by the popular jestbooks, by the *Palace of Pleasure* (an anthology of stories, mostly European), and by the few rough stories written by Englishmen (e.g., Barnabe Rich's *Farewell to Military Profession*, 1581).

During the Middle Ages short fiction had become primarily an amusing and diverting medium. The Renaissance and Enlightenment, however, made different demands of the form. The awakening concern with secular issues called for a new attention to actual conditions. Simply, the diverting stories were no longer relevant or viable. At first only the journalists and pamphleteers responded to the new demand. Short fiction disappeared, in effect, because it did not respond. When it did shake off its escapist trappings in the 19th century, it reappeared as the "modern short story." This was a new stage in the evolution of short fiction, one in which the short form undertook a new seriousness and gained a new vitality and respect.

Emergence of the modern short story. *The 19th century.* The modern short story emerged almost simultaneously in Germany, the United States, France, and Russia. In Germany there had been relatively little difference between the stories of the late 18th century and those in the older tradition of Boccaccio. In 1795 Goethe contributed a set of stories to Schiller's journal, *Die Horen*, that were obviously created with the *Decameron* in mind. Significantly, Goethe did not call them "short stories" (*Novellen*) although the term was available to him. Rather, he thought of them as "entertainments" for German travellers (*Unterhaltungen deutscher Ausgewanderten*). Friedrich Schlegel's early discussion of the short narrative form, appearing soon after Goethe's "entertainments," also focussed on Boccaccio (*Nachrichten von den poetischen Werken des G. Boccaccio*, 1801).

But a new type of short fiction was near at hand—a type that accepted some of the realistic properties of popular journalism. In 1827, 32 years after publishing his own

"entertainments," Goethe commented on the difference between the newly emergent story and the older kind. "What is a short story," he asked, "but an event which, though unheard of, has occurred? Many a work which passes in Germany under the title 'short story' is not a short story at all, but merely a tale or what else you would like to call it." Two influential critics, Christoph Wieland and Friedrich Schleiermacher, also argued that a short story properly concerned itself with events that actually happened or could happen. A short story, for them, had to be realistic.

Perhaps sensitive to this qualification, Heinrich von Kleist and E.T.A. Hoffmann called their short works on fabulous themes "tales" (*Erzählungen*). Somewhat like Poe, Kleist created an expression of human problems, partly metaphysical and partly psychological, by dramatizing man's confrontations with a fantastic, chaotic world. Hoffmann's intriguing tales of exotic places and of supernatural phenomena were very likely his most influential. Another important writer, Ludwig Tieck, explicitly rejected realism as the definitive element in a short story. As he noted in his preface to the 1829 collection of his works and as he demonstrated in his stories, Tieck envisioned the short story as primarily a matter of intensity and ironic inversion. A story did not have to be realistic in any outward sense, he claimed, so long as the chain of consequences was "entirely in keeping with character and circumstances." By allowing the writer to pursue an inner, and perhaps bizarre, reality and order, Tieck and the others kept the modern story open to nonjournalistic techniques.

In the United States, the short story, as in Germany, evolved in two strains. On the one hand there appeared the realistic story that sought objectively to deal with seemingly real places, events, or persons. The regionalist stories of the second half of the 19th century (including those by G.W. Cable, Bret Harte, Sarah Orne Jewett) are of this kind. On the other hand, there developed the impressionist story, a tale shaped and given meaning by the consciousness and psychological attitudes of the narrator. Predicated upon this element of subjectivity, these stories seem less objective and are less realistic in the outward sense. Of this sort are Poe's tales in which the hallucinations of a central character or narrator provide the details and facts of the story. Like the narrators in "The Tell-Tale Heart" (1843) and "The Imp of the Perverse" (1845), the narrator of "The Fall of the House of Usher" (1839) so distorts and transforms what he sees that the reader cannot hope to look objectively at the scene. Looking through an intermediary's eyes, the reader can see only the narrator's impressions of the scene.

Some writers contributed to the development of both types of story. Washington Irving wrote several realistic sketches (*The Sketch-Book*, 1819–20; *The Alhambra*, 1832) in which he carefully recorded appearances and actions. Irving also wrote stories in which the details were taken not from ostensible reality but from within a character's mind. Much of the substance of "The Stout Gentleman" (1821), for example, is reshaped and recharged by the narrator's fertile imagination; "Rip Van Winkle" (1819) draws upon the symbolic surreality of Rip's dreams.

The short prose of Nathaniel Hawthorne illustrates that neither type of modern story, however, has exclusive rights to the use of symbol. On a few occasions, as in "My Kinsman, Major Molineux" (1832), Hawthorne's stories are about symbolic events as they are viewed subjectively by the central character. Hawthorne's greater gift, however, was for creating scenes, persons, and events that strike the reader as being actual historical facts and also as being rich in symbolic import. "Endicott and the Red Cross" (1837) may seem little more than a photographic sketch of a tableau out of history (the 17th-century Puritan leader cuts the red cross of St. George out of the colonial flag, the first act of rebellion against England), but the details are symbols of an underground of conflicting values and ideologies.

Several American writers, from Poe to James, were interested in the "impressionist" story that focusses on the impressions registered by events on the characters' minds,

A new concern

Evolving in two strains

A new type of short fiction

The "impressionist" story

rather than the objective reality of the events themselves. In Herman Melville's "Bartleby the Scrivener" (1856) the narrator is a man who unintentionally reveals his own moral weaknesses through his telling of the story of Bartleby. Mark Twain's tales of animals ("The Celebrated Jumping Frog," 1865; "The Story of Old Ram," 1872; "Baker's Blue Jay Yarn," 1879), all impressionist stories, distort ostensible reality in a way that reflects on the men who are speaking. Ambrose Bierce's famous "An Occurrence at Owl Creek Bridge" (1891) is another example of this type of story in which the reader sees a mind at work—distorting, fabricating, and fantasizing—rather than an objective picture of actuality. In contrast, William Dean Howells usually sought an objectifying aesthetic distance. Though Howells was as interested in human psychology and behaviour as any of the impressionist writers, he did not want his details filtered through a biased, and thus distorting, narrator. Impressionism, he felt, gave license for falsifications; in the hands of many writers of his day, it did in fact result in sentimental romanticizing.

But in other hands the impressionist technique could subtly delineate human responses. Henry James was such a writer. Throughout his prefaces to the New York edition of his works, the use of an interpreting "central intelligence" is constantly emphasized. "Again and again, on review," James observes, "the shorter things in especial that I have gathered into [the Edition] have ranged themselves not as my own impersonal account of the affair in hand, but as my account of somebody's impression of it." This use of a central intelligence, who is the "impersonal author's concrete deputy or delegate" in the story, allows James all the advantages of impressionism and, simultaneously, the freedom and mobility common to stories narrated by a disembodied voice.

In at least one way, 19th-century America resembled 16th-century Italy: there was an abundance of second- and third-rate short stories. And, yet, respect for the form grew substantially, and most of the great artists of the century were actively participating in its development. The seriousness with which many writers and readers regarded the short story is perhaps most clearly evident in the amount and kind of critical attention it received. James, Howells, Harte, Twain, Melville, and Hawthorne all discussed it as an art form, usually offering valuable insights, though sometimes shedding more light on their own work than on the art as a whole.

But the foremost American critic of the short story was Edgar Allan Poe. Himself a creator of influential impressionist techniques, Poe believed that the definitive characteristic of the short story was its unity of effect. "A skillful literary artist has constructed a tale," Poe wrote in his review of Hawthorne's *Twice-Told Tales* in 1842.

If wise, he has not fashioned his thoughts to accommodate his incidents; but having conceived, with deliberate care, a certain unique or single effect to be wrought out, he then invents such incidents—he then combines such events as may best aid him in establishing this preconceived effect. If his very initial sentence tend not to the out-bringing of this effect, then he has failed in his first step. In the whole composition there should be no word written of which the tendency, direct or indirect, is not to the one pre-established design.

Poe's polemic primarily concerns craftsmanship and artistic integrity; it hardly prescribes limits on subject matter or dictates technique. As such, Poe's thesis leaves the story form open to experimentation and to growth while it demands that the form show evidence of artistic diligence and seriousness.

The new respect for the short story was also evident in France, as Henry James observed, when in 1844 Prosper Mérimée with his handful of little stories was appointed to the French Academy. As illustrated by "Columbia" (1841) or "Carmen" (1845), which gained additional fame as an opera, Mérimée's stories are masterpieces of detached and dry observation, though the subject matter itself is often emotionally charged. Nineteenth-century France produced short stories as various as 19th-century America—although the impressionist tale was generally less common in France. (It is as if, not having an outstanding impressionist storyteller themselves, the French

adopted Poe, who was being ignored by the critics in his own country.) The two major French impressionist writers were Charles Nodier, who experimented with symbolic fantasies, and Gérard de Nerval, whose collection *Les Filles du feu* (1854; "Daughters of Fire") grew out of recollections of his childhood. Artists primarily known for their work in other forms also attempted the short story—novelists like Honoré de Balzac and Gustave Flaubert and poets like Alfred de Vigny and Théophile Gautier.

One of the most interesting writers of 19th-century France is Alphonse Daudet, whose stories reflect the spectrum of interest and techniques of the entire century. His earliest and most popular stories (*Lettres de mon moulin*, 1866; "Letters from My Mill") create a romantic, picturesque fantasy; his stories of the Franco-Prussian War (*Contes du Lundi*, 1873; "Monday's Tales") are more objectively realistic, and the sociological concern of his last works betrays his increasing interest in naturalistic determinism.

The greatest French storywriter, by far, is Guy de Maupassant, a master of the objective short story. Basically, Maupassant's stories are anecdotes that capture a revealing moment in the lives of middle class citizens. This crucial moment is typically recounted in a well-plotted design, though perhaps in some stories like "Boule de suif" (1880; "Ball of Tallow") and "The Necklace" (1881) the plot is too contrived, the reversing irony too neat, and the artifice too apparent. In other stories, like "The House of Madame Tellier" (1881), Maupassant's easy and fluid prose captures the innocence and the corruption of human behaviour.

During the first two decades of the 19th century in Russia, fable writing became a fad. By all accounts the most widely read fabulist was Ivan Krylov whose stories borrowed heavily from Aesop, La Fontaine, and various Germanic sources. If Krylov's tales made short prose popular in Russia, the stories of the revered poet Aleksandr Pushkin gained serious attention for the form. Somewhat like Mérimée in France (who was one of the first to translate Pushkin, Gogol, and Turgenev into French), Pushkin cultivated a detached, rather classical style for his stories of emotional conflicts (*The Queen of Spades*, 1834). Also very popular and respected was Mikhail Lermontov's "novel," *A Hero of Our Time* (1840), which actually consists of five stories that are more or less related.

But it is Nikolay Gogol who stands at the headwaters of the Russian short story; Dostoyevsky noted that all Russian short story writers "emerged from Gogol's overcoat," a punning allusion to the master's best known story. In a manner all his own, Gogol was developing impressionist techniques in Russia simultaneously with Poe in America. Gogol published his *Arabesques* (1835) five years before Poe collected some of his tales under a similar title. Like those of Poe, Gogol's tales of hallucination, confusing reality and dream, are among his best stories ("Nevsky Prospect" and "Diary of a Madman," both 1835). The single most influential story in the first half of the 19th century in Russia was undoubtedly Gogol's "Overcoat" (1842). Blending elements of realism (natural details from the characters' daily lives) with elements of fantasy (the central character returns as a ghost), Gogol's story seems to anticipate both the impressionism of Dostoyevsky's "Underground Man" and the realism of Tolstoy's "Ivan Ilich."

Ivan Turgenev appears, at first glance, antithetical to Gogol. In *A Sportsman's Notebook* (1852) Turgenev's simple use of language, his calm pace, and his restraint clearly differentiate him from Gogol. But like Gogol, Turgenev was more interested in capturing qualities of people and places than in building elaborate plots. A remaining difference between the two Russians, however, tends to make Turgenev more acceptable to 20th-century readers: Turgenev studiously avoided anything artificial. Though he may have brought into his realistic scenes a tale of a ghost ("Bezhin Meadow," 1852), he did not attempt to bring in a ghost (as Gogol had done in "The Overcoat"). In effect, Turgenev's allegiance was wholly to detached observation.

Developing some of the interests of Gogol, Fyodor Dostoyevsky experimented with the impressionist story. The early story "White Nights" (1848), for example, is a "Tale

A master of the objective story

Russian writers

Poe's notion of the unity of effect

French short stories

of Love from the Reminiscence of a Dreamer" as the subtitle states; the title of one of his last stories, "The Dream of the Ridiculous Man" (1877), also echoes Poe and Gogol. Though sharing Dostoyevsky's interest in human motives, Leo Tolstoy used vastly different techniques. He usually sought psychological veracity through a more detached and, presumably, objective narrator ("The Death of Ivan Ilich," 1886; "The Kreutzer Sonata," 1891). Perhaps somewhat perplexed by Tolstoy's nonimpressionist means of capturing and delineating psychological impressions, Henry James pronounced Tolstoy the masterhand of the disconnection of method from matter.

The Russian master of the objective story was Anton Chekhov. No other storywriter so consistently as Chekhov turned out first-rate works. Though often compared to Maupassant, Chekhov is much less interested in constructing a well-plotted story; nothing much actually happens in Chekhov's stories, though much is revealed about his characters and the quality of their lives. While Maupassant focusses on event, Chekhov keeps his eye on character. Stories like "The Grasshopper" (1892), "The Darling" (1898), and "In the Ravine" (1900)—to name only three—all reveal Chekhov's perception, his compassion, and his subtle humour and irony. One critic says of Chekhov that he is no moralist—he simply says "you live badly, ladies and gentlemen," but his smile has the indulgence of a very wise man.

THE 20TH CENTURY

In the first half of the 20th century the appeal of the short story continued to grow. Literally hundreds of writers—including, as it seems, nearly every major dramatist, poet, and novelist—published thousands of excellent stories. William Faulkner suggested that writers often try their hand at poetry, find it too difficult, go on to the next most demanding form, the short story, fail at that, and only then settle for the novel. In the 20th century Germany, France, Russia, and the U.S. lost what had once appeared to be their exclusive domination of the form. Innovative and commanding writers emerged in countries that had previously exerted little influence on the genre: Sicily, for example, produced Luigi Pirandello; Czechoslovakia, Franz Kafka; Japan, Akutagawa Ryūnosuke; Argentina, Jorge Luis Borges. Literary journals with international circulation, such as Ford Madox Ford's *Transatlantic Review*, *Scribner's Magazine*, and Harriet Weaver's *Egoist*, provided a steady and prime exposure for young writers.

Increasing complexity of the short story

As the familiarity with it increased, the short story form itself became more varied and complex. The fundamental means of structuring a story underwent a significant change. The overwhelming or unique event that usually informed the 19th-century story fell out of favour with the storywriter of the early 20th century. He grew more interested in subtle actions and unspectacular events. Sherwood Anderson, one of the most influential U.S. writers of the early 20th century, observed that the common belief in his day was that stories had to be built around a plot, a notion that, in Anderson's opinion, appeared to poison all storytelling. His own aim was to achieve form, not plot, although form was more elusive and difficult. The record of the short story in the 20th

century is dominated by this increased sensitivity to—and experimentation with—form. Although the popular writers of the century (like O. Henry in the U.S. and Paul Morand in France) may have continued to structure stories according to plot, the greater artists turned elsewhere for structure, frequently eliciting the response from cursory readers that "nothing happens in these stories." Narratives like Ernest Hemingway's "A Clean Well-Lighted Place" may seem to have no structure at all, so little physical action develops; but stories of this kind are actually structured around a psychological, rather than physical, conflict. In several of Hemingway's stories (as in many by D.H. Lawrence, Katherine Mansfield, and others), physical action and event are unimportant except insofar as the actions reveal the psychological underpinnings of the story. Stories came to be structured, also, in accordance with an underlying archetypal model: the specific plot and characters are important insofar as they allude to a traditional plot or figure, or to patterns that have recurred with wide implications in the history of mankind. Katherine Anne Porter's "Flowering Judas," for example, echoes and ironically inverts the traditional Christian legend. Still other stories are formed by means of motif, usually a thematic repetition of an image or detail that represents the dominant idea of the story. "The Dead," the final story in James Joyce's *Dubliners*, builds from a casual mention of death and snow early in the story to a culminating paragraph that links them in a profound vision. Seldom, of course, is the specific structure of one story appropriate for a different story. Faulkner, for example, used the traditional pattern of the knightly quest (in an ironic way) for his story "Was," but for "Barn Burning" he relied on a psychologically organic form to reveal the story of young Sarty Snopes.

No single form provided the 20th-century writer with the answer to structural problems. As the primary structuring agent, spectacular and suspenseful action was rather universally rejected around midcentury since motion pictures and television could present it much more vividly. As the periodicals that had supplied escapist stories to mass audiences declined, the short story became the favoured form of a smaller but intellectually more demanding readership. The Argentine Borges, for example, attracted an international following with his *Ficciones*, stories that involved the reader in dazzling displays of erudition and imagination, unlike anything previously encountered in the genre. Similarly, the American Donald Barthelme's composition consisted of bits and pieces of, e.g., television commercials, political speeches, literary allusions, eavesdropped conversations, graphic symbols, dialogue from Hollywood movies—all interspersed with his own original prose in a manner that defied easy comprehension and yet compelled the full attention of the reader. The short story also lent itself to the rhetoric of student protest in the 1960s and was found in a bewildering variety of mixed-media forms in the "underground" press that publicized this life style throughout the world. In his deep concern with such a fundamental matter as form, the 20th-century writer unwittingly affirmed the maturation and popularity of the genre; only a secure and valued (not to mention flexible) genre could withstand and, moreover, encourage such experimentation. (A.J.H.)

The decline of the plot

DRAMA

Dramatic literature

The term dramatic literature implies a contradiction in that "literature" originally meant something written and "drama" meant something performed. Most of the problems, and much of the interest, in the study of dramatic literature stem from this contradiction. Even though a play may be appreciated solely for its qualities as writing, greater rewards probably accrue to those who remain alert to the volatility of the play as a whole.

In order to appreciate this complexity in drama, however, each of its elements—acting, directing, staging, etc.—

should be studied, so that its relationship to all the others can be fully understood. It is the purpose of this section to study drama with particular attention to what the playwright sets down. A similar approach is taken in the following sections on the two main types of dramatic literature, *Tragedy*, and *Comedy*. The history of dramatic literature is discussed in such articles as THEATRE, THE HISTORY OF WESTERN; and EAST ASIAN ARTS; as well as in articles on the history of literature such as LITERATURE, THE HISTORY OF WESTERN. Regional studies of both the theatre and literature will be found in articles such as SOUTH ASIAN ARTS; and AFRICAN ARTS.

GENERAL CHARACTERISTICS

Essential
elements of
a play

From the inception of a play in the mind of its author to the image of it that an audience takes away from the theatre, many hands and many physical elements help to bring it to life. Questions therefore arise as to what is and what is not essential to it. Is a play what its author thought he was writing, or the words he wrote? Is a play the way in which those words are intended to be embodied, or their actual interpretation by a director and his actors on a particular stage? Is a play in part the expectation an audience brings to the theatre, or is it the real response to what is seen and heard? Since drama is such a complex process of communication, its study and evaluation is as uncertain as it is mercurial.

All plays depend upon a general agreement by all participants—author, actors, and audience—to accept the operation of theatre and the conventions associated with it, just as players and spectators accept the rules of a game. Drama is a decidedly unreal activity, which can be indulged only if everyone involved admits it. Here lies some of the fascination of its study. For one test of great drama is how far it can take the spectator beyond his own immediate reality and to what use this imaginative release can be put. But the student of drama must know the rules with which the players began the game before he can make this kind of judgment. These rules may be conventions of writing, acting, or audience expectation. Only when all conventions are working together smoothly in synthesis, and the make-believe of the experience is enjoyed passionately with mind and emotion, can great drama be seen for what it is: the combined work of a good playwright, good players, and a good audience who have come together in the best possible physical circumstances.

Drama in some form is found in almost every society, primitive and civilized, and has served a wide variety of functions in the community. There are, for example, records of a sacred drama in Egypt 2,000 years before Christ, and Thespis in the 6th century BC in ancient Greece is accorded the distinction of being the first known playwright. Elements of drama such as mime and dance, costume and decor long preceded the introduction of words and the literary sophistication now associated with a play. Moreover, such basic elements were not superseded by words, merely enhanced by them. Nevertheless, it is only when a playscript assumes a disciplinary control over the dramatic experience that the student of drama gains measurable evidence of what was intended to constitute the play. Only then can dramatic literature be discussed as such.

The late
arrival of
words in
drama

The texts of plays indicate the different functions they served at different times. Some plays embraced nearly the whole community in a specifically religious celebration, as when all the male citizens of a Greek city-state came together to honour their gods; or when the annual Feast of Corpus Christi was celebrated with the great medieval Christian mystery cycles. On the other hand, the ceremonious temple ritual of the early Nō drama of Japan was performed at religious festivals only for the feudal aristocracy. But the drama may also serve a more directly didactic purpose, as did the morality plays of the later Middle Ages, some 19th-century melodramas, and the 20th-century discussion plays of George Bernard Shaw and Bertolt Brecht. Plays can satirize society, or they can gently illuminate human weakness; they can divine the greatness and the limitations of man in tragedy, or, in modern naturalistic playwriting, probe his mind. Drama is the most wide-ranging of all the arts: it not only represents life but also is a way of seeing it. And it repeatedly proves Dr. Samuel Johnson's contention that there can be no certain limit to the modes of composition open to the dramatist.

Common elements of drama. Despite the immense diversity of drama as a cultural activity, all plays have certain elements in common. For one thing, drama can never become a "private" statement—in the way a novel or a poem may be—without ceasing to be meaningful theatre. The characters may be superhuman and godlike in appearance, speech, and deed or grotesque and ridiculous, perhaps even puppets, but as long as they behave in even

vaguely recognizable human ways the spectator can understand them. Only if they are too abstract do they cease to communicate as theatre. Thus, the figure of Death in medieval drama reasons like a human being, and a god in Greek tragedy or in Shakespeare talks like any mortal. A play, therefore, tells its tale by the imitation of human behaviour. The remoteness or nearness of that behaviour to the real life of the audience can importantly affect the response of that audience: it may be in awe of what it sees, or it may laugh with detached superiority at clownish antics, or it may feel sympathy. These differences of alienation or empathy are important, because it is by opening or closing this aesthetic gap between the stage and the audience that a dramatist is able to control the spectator's experience of the play and give it purpose.

The second essential is implicit in the first. Although static figures may be as meaningfully symbolic on a stage as in a painting, the deeper revelation of character, as well as the all-important control of the audience's responses, depends upon a dynamic presentation of the figures in action. A situation must be represented on the stage, one recognizable and believable to a degree, which will animate the figures as it would in life. Some argue that action is the primary factor in drama, and that character cannot emerge without it. Since no play exists without a situation, it appears impossible to detach the idea of a character from the situation in which he is placed, though it may seem possible after the experience of the whole play. Whether the playwright conceives character before situation, or vice versa, is arbitrary. More relevant are the scope and scale of the character-in-situation—whether, for example, it is man confronting God or man confronting his wife—for that comes closer to the kind of experience the play is offering its audience. Even here one must beware of passing hasty judgment, for it may be that the grandest design for heroic tragedy may be less affecting than the teasing vision of human madness portrayed in a good farce.

A third factor is style. Every play prescribes its own style, though it will be influenced by the traditions of its theatre and the physical conditions of performance. Style is not something imposed by actors upon the text after it is written, nor is it superficial to the business of the play. Rather, it is self-evident that a play will not communicate without it. Indeed, many a successful play has style and little else. By "style," therefore, is implied the whole mood and spirit of the play, its degree of fantasy or realism, its quality of ritualism or illusion, and the way in which these qualities are signalled by the directions, explicit or implicit, in the text of the play. In its finer detail, a play's style controls the kind of gesture and movement of the actor, as well as his tone of speech, its pace and inflexion. In this way the attitude of the audience is prepared also: nothing is more disconcerting than to be misled into expecting either a comedy or a tragedy and to find the opposite, although some great plays deliberately introduce elements of both. By means of signals of style, the audience may be led to expect that the play will follow known paths, and the pattern of the play will regularly echo the rhythm of response in the auditorium. Drama is a conventional game, and spectators cannot participate if the rules are constantly broken.

The role of
style

By presenting animate characters in a situation with a certain style and according to a given pattern, a playwright will endeavour to communicate his thoughts and feelings and have his audience consider his ideas or reproduce the emotion that drove him to write as he did. In theatrical communication, however, audiences remain living and independent participants. In the process of performance, an actor has the duty of interpreting his author for the people watching him, and will expect to receive "feedback" in turn. The author must reckon with this in his writing. Ideas will not be accepted, perhaps, if they are offered forthrightly; and great dramatists who are intent on furthering social or political ideas, such as Henrik Ibsen, George Bernard Shaw, and Bertolt Brecht, quickly learned methods of having the spectator reason the ideas for himself as part of his response to the play. Nor will passions necessarily be aroused if overstatement of feeling

("sentimentality") is used without a due balance of thinking and even the detachment of laughter: Shakespeare and Chekhov are two outstanding examples in Western drama of writers who achieved an exquisite balance of pathos with comedy in order to ensure the affective function of their plays.

Dramatic expression. The language of drama can range between great extremes: on the one hand, an intensely theatrical and ritualistic manner; and on the other, an almost exact reproduction of real life of the kind commonly associated with motion picture and television drama. In the ritualistic drama of ancient Greece, the playwrights wrote in verse, and it may be assumed that their actors rendered this in an incantatory speech halfway between speech and song. Both the popular and the coterie drama of the Chinese and Japanese theatre were also essentially operatic, with a lyrical dialogue accompanied by music and chanted rhythmically. The effect of such rhythmical delivery of the words was to lift the mood of the whole theatre onto the level of religious worship. Verse is employed in other drama that is conventionally elevated, like the Christian drama of the Middle Ages, the tragedy of the English Renaissance, the heroic Neoclassical tragedies of 17th-century France by Pierre Corneille and Jean Racine, the Romantic lyricism of Goethe and Schiller, and modern attempts at a revival of a religious theatre like those of T.S. Eliot. Indeed, plays written in prose dialogue were at one time comparatively rare, and then associated essentially with the comic stage. Only at the end of the 19th century, when naturalistic realism became the mode, were characters in dramas expected to speak as well as behave as in real life.

Elevation is not the whole rationale behind the use of verse in drama. Some critics maintain that a playwright can exercise better control both over the speech and movement of his actors and over the responses of his audience by using the more subtle tones and rhythms of good poetry. The loose, idiomatic rhythms of ordinary conversation, it is argued, give both actor and spectator too much freedom of interpretation and response. Certainly, the aural, kinetic, and emotive directives in verse are more direct than prose, though, in the hands of a master of prose dialogue like Shaw or Chekhov, prose can also share these qualities. Even more certain, the "aesthetic distance" of the stage, or the degree of unreality and make-believe required to release the imagination, is considerably assisted if the play uses elements of verse, like rhythm and rhyme, not found in ordinary speech. Thus, verse drama may embrace a wide variety of nonrealistic aural and visual devices: Greek tragic choric speech provided a philosophical commentary upon the action, which at the same time drew the audience lyrically into the mood of the play. In the drama of India, a verse accompaniment made the actors' highly stylized system of symbolic gestures of head and eyes, arms and fingers a harmonious whole. The tragic soliloquy in Shakespeare permitted the hero, alone on the stage with his audience, to review his thoughts aloud in the persuasive terms of poetry; thus, the soliloquy was not a stopping place in the action but rather an engrossing moment of drama when the spectator's mind could leap forward.

Dramatic structure. The elements of a play do not combine naturally to create a dramatic experience but, rather, are made to work together through the structure of a play, a major factor in the total impact of the experience. A playwright will determine the shape of a play in part according to the conditions in which it will be performed: how long should it take to engage an audience's interest and sustain it? How long can an audience remain in their seats? Is the audience sitting in one place for the duration of performance, or is it moving from one pageant stage to the next, as in some medieval festivals? Structure is also dictated by the particular demands of the material to be dramatized: a revue sketch that turns on a single joke will differ in shape from a religious cycle, which may portray the whole history of mankind from the Creation to the Last Judgment. A realistic drama may require a good deal of exposition of the backgrounds and memories of the characters, while in a chronicle play the playwright

may tell the whole story episodically from its beginning to the end. There is one general rule, as Aristotle originally suggested in his *Poetics*: a play must be long enough to supply the information an audience needs to be interested and to generate the experience of tragedy, or comedy, on the senses and imagination.

In the majority of plays it is necessary to establish a conventional code of place and time. In a play in which the stage must closely approximate reality, the location of the action will be precisely identified, and the scenic representation on stage must confirm the illusion. In such a play, stage time will follow chronological time almost exactly; and if the drama is broken into three, four, or five acts, the spectator will expect each change of scene to adjust the clock or the calendar. But the theatre has rarely expected realism, and by its nature it allows an extraordinary freedom to the playwright in symbolizing location and duration: as Dr. Samuel Johnson observed in his discussion of this freedom in Shakespeare, the spectators always allow the play to manipulate the imagination. It is sufficient for the witches in *Macbeth* to remark their "heath" with its "fog and filthy air" for their location to be accepted on a stage without scenery; and when Lady Macbeth later is seen alone reading a letter, she is without hesitation understood to be in surroundings appropriate to the wife of a Scottish nobleman. Simple stage symbolism may assist the imagination, whether the altar of the gods situated in the centre of the Greek *orchestra*, a strip of red cloth to represent the Red Sea in a medieval miracle play, or a chair on which the Tibetan performer stands to represent a mountain. With this degree of fantasy, it is no wonder that the theatre can manipulate time as freely, passing from the past to the future, from this world to the next, and from reality to dream.

It is questionable, therefore, whether the notion of "action" in a play describes what happens on the stage or what is recreated in the mind of the audience. Certainly it has little to do with merely physical activity by the players. Rather, anything that urges forward the audience's image of the play and encourages the growth of its imagination is a valid part of the play's action. Thus, it was sufficient for the ancient Greek dramatist Aeschylus to have only two speaking male actors who wore various masks, typed for sex, age, class, and facial expression. In the Italian 16th- and 17th-century *commedia dell'arte*, the standard characters Pantalone and Arlecchino, each wearing his traditional costume and mask, appeared in play after play and were immediately recognized, so that an audience could anticipate the behaviour of the grasping old merchant and his rascally servant. On a less obvious level, a speech that in reading seems to contribute nothing to the action of the play can provide in performance a striking stimulus to the audience's sense of the action, its direction and meaning. Thus, both the Greek chorus and the Elizabethan actor in soliloquy might be seen to "do" nothing, but their intimate speeches of evaluation and reassessment teach the spectator how to think and feel about the action of the main stage and lend great weight to the events of the play. For drama is a reactive art, moving constantly in time, and any convention that promotes a deep response while conserving precious time is of immeasurable value.

DRAMA AS AN EXPRESSION OF A CULTURE

In spite of the wide divergencies in purpose and convention of plays as diverse as the popular kabuki of Japan and the coterie comedies of the Restoration in England, a Javanese puppet play and a modern social drama by the contemporary American dramatist Arthur Miller, all forms of dramatic literature have some points in common. Differences between plays arise from differences in conditions of performance, in local conventions, in the purpose of theatre within the community, and in cultural history. Of these, the cultural background is the most important, if the most elusive. It is cultural difference that makes the drama of the East immediately distinguishable from that of the West.

East-West differences. Oriental drama consists chiefly of the classical theatre of Hindu India and its derivatives in Malaya and of Burma, Thailand, China, Japan, Java,

The role of
verse

Conven-
tions of
time and
place

Conven-
tions of
action

and Bali. It was at its peak during the period known in the West as the Middle Ages and the Renaissance. Stable and conservative, perpetuating its customs with reverence, Oriental culture showed little of the interest in chronology and advancement shown by the West and placed little emphasis on authors and their individual achievements. Thus the origins of the drama of the Orient are lost in time, although its themes and characteristic styles probably remain much the same as before records were kept. The slow-paced, self-contained civilizations of the East have only recently been affected by Western theatre, just as the West has only recently become conscious of the theatrical wealth of the East and what it could do to fertilize the modern theatre (as in the 20th-century experimental drama of William Butler Yeats and Thornton Wilder in English, of Paul Claudel and Antonin Artaud in French, and of Bertolt Brecht in German).

Stylization
of Oriental
theatre

In its representation of life, classical Oriental drama is the most conventional and nonrealistic in world theatre. Performed over the centuries by actors devoted selflessly to the profession of a traditional art, conventions of performance became highly stylized, and traditions of characterization and play structure became formalized to a point of exceptional finesse, subtlety, and sophistication. In Oriental drama all the elements of the performing arts are made by usage to combine to perfection: dance and mime, speech and song, narrative and poetry. The display and studied gestures of the actors, their refined dance patterns, and the all-pervasive instrumental accompaniment to the voices of the players and the action of the play, suggest to Western eyes an exquisite combination of ballet with opera, in which the written text assumes a subordinate role. In this drama, place could be shifted with a license that would have astonished the most romantic of Elizabethan dramatists, the action could leap back in time in a way reminiscent of the "flashback" of the modern cinema, and events could be telescoped with the abandon of modern expressionism. This extreme theatricality lent an imaginative freedom to its artists and audiences upon which great theatre could thrive. Significantly, most Oriental cultures also nourished a puppet theatre, in which stylization of character, action, and staging were particularly suitable to marionettes. In the classical puppet theatre of Japan, the *bunraku*, the elocutionary art of a chanted narration and the manipulative skill with the dolls diminished the emphasis on the script except in the work of the 17th-century master Chikamatsu, who enjoyed a creative freedom in writing for puppets rather than for the actors of the Kabuki. By contrast, Western drama during and after the Renaissance has offered increasing realism, not only in decor and costume but also in the treatment of character and situation.

It is generally thought that Oriental drama, like that of the West, had its beginnings in religious festivals. Dramatists retained the moral tone of religious drama while using popular legendary stories to imbue their plays with a romantic and sometimes sensational quality. This was never the sensationalism of novelty that Western dramatists sometimes used: Eastern invention is merely a variation on what is already familiar, so that the slightest changes of emphasis could give pleasure to the cognoscenti. This kind of subtlety is not unlike that found in the repeatedly depicted myths of Greek tragedy. What is always missing in Oriental drama is that restlessness for change characteristic of modern Western drama. In the West, religious questioning, spiritual disunity, and a belief in the individual vision combined finally with commercial pressures to produce comparatively rapid changes. None of the moral probing of Greek tragedy, the character psychology of Shakespeare and Racine, the social and spiritual criticism of Ibsen and Strindberg, nor the contemporary drama of shock and argument, is imaginable in the classical drama of the East.

Drama in Western cultures. The form and style of ancient Greek tragedy, which flowered in the 5th century BC in Athens, was dictated by its ritual origins and by its performance in the great dramatic competitions of the spring and winter festivals of Dionysus. Participation in ritual requires that the audience largely knows what to expect.

Ritual dramas were written on the same legendary stories of Greek heroes in festival after festival. Each new drama provided the spectators with a reassessment of the meaning of the legend along with a corporate religious exercise. Thus, the chorus of Greek tragedy played an important part in conveying the dramatist's intention. The chorus not only provided a commentary on the action but also guided the moral and religious thought and emotion of the audience throughout the play: for Aeschylus (c. 525–456 BC) and Sophocles (c. 496–406 BC) it might be said that the chorus *was* the play, and even for Euripides (c. 480–406 BC) it remained lyrically powerful. Other elements of performance also controlled the dramatist in the form and style he could use in these plays: in particular, the great size of the Greek arena demanded that the players make grand but simple gestures and intone a poetry that could never approach modern conversational dialogue. Today, the superhuman characters of these plays, Agamemnon and Clytemnestra, Orestes and Electra, Oedipus and Antigone, seem unreal, for they display little "characterization" in the modern sense and their fates are sealed. Nevertheless, these great operatic tableaux, built, as one critic has said, for weight and not speed, were evidently able to carry their huge audiences to a catharsis of feeling. It is a mark of the piety of those audiences that the same reverent festivals supported a leavening of satyr-plays and comedies, bawdy and irreverent comments on the themes of the tragedies, culminating in the wildly inventive satires of Aristophanes (c. 445–c. 385 BC.)

The study of Greek drama demonstrates how the ritual function of theatre shapes both play and performance. This ritual aspect was lost when the Romans assimilated Greek tragedy and comedy. The Roman comedies of Plautus (c. 254–184 BC) and Terence (c. 186/185–159 BC) were brilliant but inoffensive entertainments, while the oratorical tragedies of Seneca (c. 4 BC–AD 65) on themes from the Greek were written probably only to be read by the ruling caste. Nevertheless, some of the dramatic techniques of these playwrights influenced the shape and content of plays of later times. The bold prototype characters of Plautus (the boasting soldier, the old miser, the rascally parasite), with the intricacies of his farcical plotting, and the sensational content and stoical attitudes of Seneca's drama reappeared centuries later when classical literature was rediscovered.

Western drama had a new beginning in the medieval church, and, again, the texts reflect the ritual function of the theatre in society. The Easter liturgy, the climax of the Christian calendar, explains much of the form of medieval drama as it developed into the giant mystery cycles. From at least the 10th century the clerics of the church enacted the simple Latin liturgy of the *Quem quaeritis?* (literally "Whom do you seek?"), the account of the visit to Jesus Christ's tomb by the three Marys, who are asked this question by an angel. The liturgical form of Lent and the Passion, indeed, embodies the drama of the Resurrection to be shared mutually by actor-priest and audience-congregation. When the Feast of Corpus Christi was instituted in 1246, the great lay cycles of Biblical plays (the mystery or miracle cycles) developed rapidly, eventually treating the whole story of man from the Creation to the Last Judgment, with the Crucifixion still the climax of the experience. The other influence controlling their form and style was their manner of performance. The vast quantity of material that made up the story was broken into many short plays, and each was played on its own stage in the vernacular by members of the craft guilds. Thus, the authors of these dramas gave their audience not a mass communal experience, as the Greek dramatists had done, but rather many small and intimate dramatizations of the Bible story. In stylized and alliterative poetry, they mixed awesome events with moments of extraordinary simplicity, embodying local details, familiar touches of behaviour, and the comedy and the cruelty of medieval life. Their drama consists of strong and broad contrasts, huge in perspective but meaningful in human terms, religious and appropriately didactic in content and yet popular in its manner of reaching its simple audiences.

In an account of dramatic literature, the ebullient but

Use of
familiar
legends

Establishment
of prototypes
of characters

The influence of improvisation

unscripted farces and romances of the *commedia dell'arte* properly have no place, but much in it became the basis of succeeding comedy. Two elements are worth noting. First, the improvisational spirit of the *commedia* troupes, in which the actor would invent words and comic business (*lazzi*) to meet the occasion of the play and the audience he faced, encouraged a spontaneity in the action that has affected the writing and playing of Western comedy ever since. Second, basic types of comic character derived from the central characters, who reappeared in the same masks in play after play. As these characters became well known everywhere, dramatists could rely on their audience to respond to them in predictable fashion. Their masks stylized the whole play and allowed the spectator freedom to laugh at the unreality of the action. An understanding of the *commedia* illuminates a great deal in the written comedies of Shakespeare in England, of Molière and Marivaux in France, and of Goldoni and Gozzi in Italy.

In the 16th century, England and Spain provided all the conditions necessary for a drama that could rival ancient Greek drama in scope and subtlety. In both nations, there were public as well as private playhouses, audiences of avid imagination, a developing language that invited its poetic expansion, a rapid growth of professional acting companies, and a simple but flexible stage. All these factors combined to provide the dramatist with an opportunity to create a varied and exploratory new drama of outstanding interest. In Elizabethan London, dramatists wrote in an extraordinary range of dramatic genres, from native comedy and farce to Senecan tragedy, from didactic morality plays to popular chronicle plays and tragicomedies, all before the advent of Shakespeare (1564–1616). Although Shakespeare developed certain genres, such as the chronicle play and the tragedy, to a high degree, Elizabethan dramatists characteristically used a medley of styles. With the exception of Ben Jonson (1572/73–1637) and a few others, playwrights mixed their ingredients without regard for classical rule. The result was a rich body of drama, exciting and experimental in character. A host of new devices were tested, mixing laughter and passion; shifting focus and perspective by slipping from verse to prose and back again; extending the use of the popular clown; exploiting the double values implicit in boy actors playing the parts of girls; exploring the role of the actor in and out of character; but, above all, developing an extraordinarily flexible dramatic poetry. These dramatists produced a visually and aurally exciting hybrid drama that could stress every subtlety of thought and feeling. It is not surprising that they selected their themes from every Renaissance problem of order and authority, of passion and reason, of good and evil and explored every comic attitude to people and society with unsurpassed vigour and vision.

French drama in the 17th century

Quite independently in Spain, dramatists embarked upon a parallel development of genres ranging from popular farce to chivalric tragedy. The hundreds of plays of Spain's greatest playwright, Lope de Vega (1562–1635), cover every subject from social satire to religion with equal exuberance. The drama of Paris of the 17th century, however, was determined by two extremes of dramatic influence. On the one hand, some playwrights developed a tragedy rigidly based in form upon Neoclassical notions of Aristotelian unity, controlled by verse that is more regular than that of the Spanish or English dramatists. On the other hand, the French theatre developed a comedy strongly reflecting the work of the itinerant troupes of the *commedia dell'arte*. The Aristotelian influence resulted in the plays of Pierre Corneille (1606–84) and Jean Racine (1639–99), tragedies of honour using classical themes, highly sophisticated theatrical instruments capable of searching deeply into character and motive, and capable of creating the powerful tension of a tightly controlled plot. The other influence produced the brilliant plays of Molière (1622–73), whose training as an actor in the masked and balletic *commedia* tradition supplied him with a perfect mode for a more sophisticated comedy. Molière's work established the norm of French comedy, bold in plotting, exquisite in style, irresistible in comic suggestion. Soon after, upon the return of Charles II to the throne of England in 1660, a revival of theatre started the English drama on a new

course. Wits such as William Wycherley (1640–1716) and William Congreve (1670–1729) wrote for the intimate playhouses of the Restoration and an unusually homogeneous coterie audience of the court circle. They developed a "comedy of manners," replete with social jokes that the actor, author, and spectator could share—a unique phase in the history of drama. These plays started a characteristic style of English domestic comedy still recognizable in London comedy today.

German dramatists of the later part of the 18th century achieved stature through a quite different type of play: Johann Wolfgang von Goethe (1749–1832), Johann Christoph Friedrich von Schiller (1759–1805), and others of the passionate, poetic Sturm und Drang ("storm and stress") movement tried to echo the more romantic tendencies in Shakespeare's plays. Dramatists of the 19th century, however, lacking the discipline of classical form, wrote derivative melodramas that varied widely in quality, often degenerating into mere sensationalism. Melodrama rapidly became the staple of the theatre across Europe and America. Bold in plotting and characterization, simple in its evangelical belief that virtue will triumph and providence always intervene, it pleased vast popular audiences and was arguably the most prolific and successful drama in the history of the theatre. Certainly, melodrama's elements of essential theatre should not be ignored by those interested in drama as a social phenomenon. At least melodramas encouraged an expansion of theatre audiences ready for the most recent phase in dramatic history.

Melodrama

The time grew ripe for a new and more adult drama at the end of the 19th century. As novelists developed greater naturalism in both content and style, dramatists too looked to new and more realistic departures: the dialectical comedies of ideas of George Bernard Shaw (1856–1950); the problem plays associated with Henrik Ibsen (1828–1906); the more lyrical social portraits of Anton Chekhov (1860–1904); the fiercely personal, social, and spiritual visions of August Strindberg (1849–1912). These dramatists began by staging the speech and behaviour of real life, in devoted detail, but became more interested in the symbolic and poetic revelation of the human condition. Where Ibsen began by modelling his tightly structured dramas of man in society upon the formula for the "well made" play, which carefully controlled the audience's interest to the final curtain, Strindberg, a generation later, developed a free psychological and religious dream play that bordered on Expressionism. As sophisticated audiences grew interested more in causes rather than in effects, the great European playwrights of the turn of the century mixed their realism increasingly with symbolism. Thus the Naturalistic movement in drama, though still not dead, had a short but vigorous life. Its leaders freed the drama of the 20th century to pursue every kind of style, and subsequent dramatists have been wildly experimental. The playwright today can adopt any dramatic mode, mixing his effects to shock the spectator into an awareness of himself, his beliefs, and his environment.

Drama in Eastern cultures. Because of its inborn conservatism, the dramatic literature of the East does not show such diversity, despite its variety of cultures and subcultures. The major features of Oriental drama may be seen in the three great classical sources of India, China, and Japan. The simplicity of the Indian stage, a platform erected for the occasion in a palace or a courtyard, like the simplicity of the Elizabethan stage, lent great freedom to the imagination of the playwright. In the plays of India's greatest playwright, Kālidāsa (probably 4th century AD), there is an exquisite refinement of detail in presentation. His delicate romantic tales leap time and place by simple suggestion and mingle courtly humour and light-hearted wit with charming sentiment and religious piety. Quite untrammelled by realism, lyrical in tone and refined in feeling, his fanciful love and adventure stories completely justify their function as pure entertainment. His plots are without the pain of reality, and his characters never descend from the ideal: such poetic drama is entirely appropriate to the Hindu aesthetic of blissful idealism in art.

Indian theatre

Some contrast may be felt between the idealistic style of the Sanskrit drama and the broader, less courtly man-

ner of the Chinese and its derivatives in Southeast Asia. These plays cover a large variety of subjects and styles, but all combine music, speech, song, and dance, as does all Oriental drama. Heroic legends, pathetic moral stories, and brilliant farces all blended spectacle and lyricism and were as acceptable to a sophisticated court audience as to a popular street audience. The most important Chinese plays stem from the Yüan dynasty (1206–1368), in which an episodic narrative is carefully structured and unified. Each scene introduces a song whose lines have a single rhyme, usually performed by one singer, with a code of symbolic gestures and intonations that has been refined to an extreme. The plays have strongly typed heroes and villains, simple plots, scenes of bold emotion, and moments of pure mime. Chinese drama avoided both the crudity of European melodrama and the esotericism of Western coterie drama.

Nō and
Kabuki

The drama of Japan may be said to embrace both. There, the exquisite artistry of gesture and mime, and the symbolism of setting and costume, took two major directions. The Nō drama, emerging from religious ritual, maintained a special refinement appropriate to its origins and its aristocratic audiences; the Kabuki (its name suggesting its composition: *ka*, "singing"; *bu*, "dancing"; *ki*, "acting") in the 17th century became Japan's popular drama. Nō theatre is reminiscent of the religious tragedy of the Greeks in the remoteness of its legendary content, in its masked heroic characters, in its limit of two actors and a chorus, and in the static, oratorical majesty of its style. The Kabuki, on the other hand, finds its material in domestic stories and in popular history, and the actors, without masks, move and speak more freely, without seeming to be realistic. The Kabuki plays are less rarefied and are often fiercely energetic and wildly emotional as befitting their presentation before a broader audience. The written text of the Nō play is highly poetic and pious in tone, compressed in its imaginative ideas, fastidious and restrained in verbal expression, and formal in its sparse plotting; the text of a Kabuki play lends plentiful opportunities for spectacle, sensation, and melodrama. In the Kabuki there can be moments of realism, but also whole episodes of mime and acrobatics; there can be moments of slapstick, but also moments of violent passion. In all, the words are subordinate to performance in the Kabuki.

Drama and communal belief. The drama that is most meaningful and pertinent to its society is that which arises from it and is not imposed upon it. The religious drama of ancient Greece, the temple drama of early India and Japan, the mystery cycles of medieval Europe, all have in common more than their religious content: when the theatre is a place of worship, its drama goes to the roots of belief in a particular community. The dramatic experience becomes a natural extension of man's life both as an individual and as a social being. The content of the mystery cycles speaks formally for the orthodox dogma of the church, thus seeming to place the plays at the centre of medieval life, like the church itself. Within such a comprehensive scheme, particular needs could be satisfied by comic or pathetic demonstration; for example, such a crucial belief as that of the Virgin Birth of Jesus was presented in the York (England) cycle of mystery plays, of the 14th–16th centuries, with a nicely balanced didacticism when Joseph wonders how a man of his age could have got Mary with child and an Angel explains what has happened; the humour reflects the simplicity of the audience and at the same time indicates the perfect faith that permitted the near-blasphemy of the joke. In the tragedies Shakespeare wrote for the Elizabethan theatre, he had the same gift of satisfying deep communal needs while meeting a whole range of individual interests present in his audience.

Didactic
drama

When the whole community shares a common heritage, patriotic drama and drama commemorating national heroes, as are seen almost universally in the Orient, is of this kind. Modern Western attempts at a religious didactic drama, or indeed at any drama of "ideas," have had to reckon with the disparate nature of the audience. Thus the impact of Ibsen's social drama both encouraged and divided the development of the theatre in the last years of the 19th century. Plays like *A Doll's House* (1879) and

Ghosts (published 1881), which challenged the sanctity of marriage and questioned the loyalty a wife owed to her husband, took their audiences by storm: some violently rejected the criticism of their cherished social beliefs, and thus such plays may be said to have failed to persuade general audiences to examine their moral position; on the other hand, there were sufficient numbers of enthusiasts (so-called Ibsenites) to stimulate a new drama of ideas. "Problem" plays appeared all over Europe and undoubtedly rejuvenated the theatre for the 20th century. Shaw's early Ibsenite plays in London, attacking a negative drawing-room comedy with themes of slum landlordism (*Widowers' Houses*, 1892) and prostitution (*Mrs. Warren's Profession*, 1902) resulted only in failure, but Shaw quickly found a comic style that was more disarming. In his attack on false patriotism (*Arms and the Man*, 1894) and the motives for middle class marriage (*Candida*, 1897), he does not affront his audiences before leading them by gentle laughter and surprise to review their own positions.

INFLUENCES ON THE DRAMATIST

The author of a play is affected, consciously or unconsciously, by the conditions under which he conceives and writes, by his social and economic status as a playwright, by his personal background, by his religious or political position, by his purpose in writing. The literary form of the play and its stylistic elements will be influenced by tradition, a received body of theory and dramatic criticism, as well as by the author's innovative energy. Auxiliary theatre arts such as music and design also have their own controlling traditions and conventions, which the playwright must respect. The size and shape of the playhouse, the nature of its stage and equipment, and the actor-audience relationship it encourages also determine the character of the writing. Not least, the audience's cultural assumptions, holy or profane, local or international, social or political, may override all else in deciding the form and content of the drama. These are large considerations that can take the student of drama into areas of sociology, politics, social history, religion, literary criticism, philosophy and aesthetics, and beyond.

The role of theory. It is difficult to assess the influence of theory since theory usually is based on existing drama, rather than drama on theory. Philosophers, critics, and dramatists have attempted both to describe what happens and to prescribe what should happen in drama, but all their theories are affected by what they have seen and read.

Western theory. In Europe, the earliest extant work of dramatic theory, the fragmentary *Poetics* of Aristotle (384–322 BC), chiefly reflecting his views on Greek tragedy and his favorite dramatist, Sophocles, is still relevant to an understanding of the elements of drama. Aristotle's elliptical way of writing, however, encouraged different ages to place their own interpretation upon his statements and to take as prescriptive what many believe to have been meant only to be descriptive. There has been endless discussion of his concepts *mimēsis* ("imitation"), the impulse behind all the arts, and *katharsis* ("purgation," "purification of emotion"), the proper end of tragedy, though these notions were conceived, in part, in answer to Plato's attack on *poiesis* (making) as an appeal to the irrational. That "character" is second in importance to "plot" is another of Aristotle's concepts that may be understood with reference to the practice of the Greeks, but not more realistic drama, in which character psychology has a dominant importance. The concept in the *Poetics* that has most affected the composition of plays in later ages has been that of the so-called unities—that is, of time, place, and action. Aristotle was evidently describing what he observed—that a typical Greek tragedy had a single plot and action that lasts one day; he made no mention at all of unity of place. Neoclassical critics of the 17th century, however, codified these discussions into rules.

Considering the inconvenience of such rules and their final unimportance, one wonders at the extent of their influence. The Renaissance desire to follow the ancients and its enthusiasm for decorum and classification may explain it in part. Happily, the other classical work recognized at this time was Horace's *Art of Poetry* (c. 24 BC), with its

"Rules" of
drama

basic precept that poetry should offer pleasure and profit and teach by pleasing, a notion that has general validity to this day. Happily, too, the popular drama, which followed the tastes of its patrons, also exerted a liberating influence. Nevertheless, discussion about the supposed need for the unities continued throughout the 17th century (culminating in the French critic Nicolas Boileau's *Art of Poetry*, originally published in 1674), particularly in France, where a master like Racine could translate the rules into a taut, intense theatrical experience. Only in Spain, where Lope de Vega published his *New Art of Writing Plays* (1609), written out of his experience with popular audiences, was a commonsense voice raised against the classical rules, particularly on behalf of the importance of comedy and its natural mixture with tragedy. In England both Sir Philip Sidney in his *Apologie for Poetry* (1595) and Ben Jonson in *Timber* (1640) merely attacked contemporary stage practice. Jonson, in certain prefaces, however, also developed a tested theory of comic characterization (the "humours") that was to affect English comedy for a hundred years. The best of Neoclassical criticism in English is John Dryden's *Of Dramatick Poesie, an Essay* (1668). Dryden approached the rules with a refreshing honesty and argued all sides of the question; thus he questioned the function of the unities and accepted Shakespeare's practice of mixing comedy and tragedy.

The lively imitation of nature came to be acknowledged as the primary business of the playwright and was confirmed by the authoritative voices of Dr. Samuel Johnson, who said in his *Preface to Shakespeare* (1765) "there is always an appeal open from criticism to nature," and the German dramatist and critic Gotthold Ephraim Lessing, who in his *Hamburgische Dramaturgie* (or *Hamburg Dramaturgy*; 1767-69) sought to accommodate Shakespeare to a new view of Aristotle. With the classical straitjacket removed, there was a release of dramatic energies in new directions. There were still local critical skirmishes, such as Jeremy Collier's attack on the "immorality and profaneness of the English stage" in 1698; Goldoni's attacks upon the already dying Italian commedia on behalf of greater realism; and Voltaire's reactionary wish to return to the unities and to rhymed verse in French tragedy, which was challenged in turn by Diderot's call for a return to nature. But the way was open for the development of the middle class *drame* and the excursions of romanticism. Victor Hugo, in his Preface to his play *Cromwell* (1827), capitalized on the new psychological romanticism of Goethe and Schiller as well as the popularity of the sentimental *drame* in France and the growing admiration for Shakespeare; Hugo advocated truth to nature and a dramatic diversity that could yoke together the sublime and the grotesque. This view of what drama should be received support from Émile Zola in the preface to his play *Thérèse Raquin* (1873), in which he argued a theory of naturalism that called for the accurate observation of people controlled by their heredity and environment. From such sources came the subsequent intellectual approach of Ibsen and Chekhov and a new freedom for such seminal innovators of the 20th century as Luigi Pirandello, with his teasing mixtures of absurdist laughter and psychological shock; Bertolt Brecht (1898-1956), deliberately breaking the illusion of the stage; and Antonin Artaud (1896-1948), advocating a theatre that should be "cruel" to its audience, employing all and any devices that lie to hand. The modern dramatist may be grateful that he is no longer hidebound by theory and yet also regret, paradoxically, that the theatre of his time lacks those artificial limits within which an artifact of more certain efficiency can be wrought.

Eastern theory. The Oriental theatre has always had such limits, but with neither the body of theory nor the pattern of rebellion and reaction found in the West. The Sanskrit drama of India, however, throughout its recorded existence has had the supreme authority of the *Nāṭya-śāstra*, ascribed to Bharata (c. 1st century AD), an exhaustive compendium of rules for all the performing arts, but particularly for the sacred art of drama with its auxiliary arts of dance and music. Not only does the *Nāṭya-śāstra* identify many varieties of gesture and movement but it also describes the multiple patterns that drama can as-

sume, similar to a modern treatise on musical form. Every conceivable aspect of a play is treated, from the choice of metre in poetry to the range of moods a play can achieve; but perhaps its primary importance lies in its justification of the aesthetic of Indian drama as a vehicle of religious enlightenment.

In Japan, the most celebrated of early Nō writers, Zeami Motokiyo (1363-1443), left an influential collection of essays and notes to his son about his practice, and his deep knowledge of Zen Buddhism infused the Nō drama with ideals for the art that have persisted. Religious serenity of mind (*yūgen*), conveyed through an exquisite elegance in a performance of high seriousness, is at the heart of Zeami's theory of dramatic art. Three centuries later, the outstanding dramatist Chikamatsu (1653-1725) built equally substantial foundations for the Japanese puppet theatre, later known as the *bunraku*. His heroic plays for this theatre established an unassailable dramatic tradition of depicting an idealized life inspired by a rigid code of honour and expressed with extravagant ceremony and fervent lyricism. At the same time, in another vein, his pathetic "domestic" plays of middle class life and the suicides of lovers established a comparatively realistic mode for Japanese drama, which strikingly extended the range of both the *bunraku* and the Kabuki. Today, these forms, together with the more aristocratic and intellectual Nō, constitute a classical theatre based on practice rather than on theory. They may be superseded as a result of the recent invasion of Western drama, but in their perfection they are unlikely to change. The Yüan drama of China was similarly based upon a slowly evolved body of laws and conventions derived from practice, for, like the Kabuki of Japan, this too was essentially an actors' theatre, and practice rather than theory accounts for its development.

The role of music and dance. The Sanskrit treatise *Nāṭya-śāstra* suggests that drama had its origin in the art of dance, and any survey of Western theatre, too, must recognize a comparable debt to music in the classical Greek drama, which is believed to have sprung from celebratory singing to Dionysus. Similarly, the drama of the medieval church began with the chanted liturgies of the Roman mass. In the professional playhouses of the Renaissance and after, only rarely is music absent: Shakespeare's plays, particularly the comedies, are rich with song, and the skill with which he pursues dramatic ends with musical help is a study in itself. Molière conceived most of his plays as comedy-ballets, and much of his verbal style derives directly from the balletic qualities of the commedia. The popularity of opera in the 18th century led variously to John Gay's prototype for satirical ballad-opera, *The Beggar's Opera* (1728), the opera buffa in Italy, and the opéra comique in France. The development of these forms, however, resulted in the belittling of the written drama, with the notable exception of the parodistic wit of W.S. Gilbert (1836-1911). It is worth noting, however, that the most successful modern "musicals" lean heavily on their literary sources. Today, two of the strongest influences on contemporary theatre are those of Bertolt Brecht, who believed that a dialectical theatre should employ music not merely as a background embellishment but as an equal voice with the actor's, and of Antonin Artaud, who argued that the theatre experience should subordinate the literary text to mime, music, and spectacle. Since it is evident that drama often involves a balance of the arts, an understanding of their interrelationships is proper to a study of dramatic literature.

The influence of theatre design. Though apparently an elementary matter, the shape of the stage and auditorium probably offers the greatest single control over the text of the play that can be measured and tested. Moreover, it is arguable that the playhouse architecture dictates more than any other single factor the style of a play, the conventions of its acting, and the quality of dramatic effect felt by its audience. The shape of the theatre is always changing, so that to investigate its function is both to understand the past and to anticipate the future. Today, Western theatre is in the process of breaking away from the dominance of the Victorian picture-frame theatre, and therefore from the kind of experience this produced.

Three basic
playhouse
types

The contemporary English critic John Wain has called the difference between Victorian and Elizabethan theatre a difference between "consumer" and "participation" art. The difference resulted from the physical relationship between the audience and the actor in the two periods, a relationship that determined the kind of communication open to the playwright and the role the drama could play in society. Three basic playhouse shapes have emerged in the history of the theatre: the arena stage, the open stage, and the picture-frame.

The arena stage. To the arena, or theatre-in-the-round, belongs the excitement of the circus, the bullring, and such sports as boxing and wrestling. Arena performance was the basis for all early forms of theatre—the Druid ceremonies at Stonehenge, the Tibetan harvest-festival drama, probably early Greek ritual dancing in the *orchestra*, the medieval rounds in 14th-century England and France, the medieval street plays on pageant wagons, the early Nō drama of Japan, the royal theatre of Cambodia. Characteristic of all these theatres is the bringing together of whole communities for a ritual experience; therefore, a sense of ritualistic intimacy and involvement is common to the content of the drama, and only the size of the audience changes the scale of the sung or spoken poetry. Clearly, the idiom of realistic dialogue would have been inappropriate both to the occasion and the manner of such theatre.

The open stage. When more narrative forms of action appeared in drama and particular singers or speakers needed to control the attention of their audience by facing them, the open, "thrust," or platform stage, with the audience on three sides of the actor, quickly developed its versatility. Intimate and ritualistic qualities in the drama could be combined with a new focus on the players as individual characters. The open stage and its variants were used by the majority of great national theatres, particularly those of China and Japan, the booths of the Italian commedia, the Elizabethan public and private playhouses, and the Spanish *corrales* (i.e., the areas between town houses) of the Renaissance. While open-stage performance discouraged scenic elaboration, it stressed the actor and his role, his playing to and away from the spectators, with the consequent subtleties of empathy and alienation. It permitted high style in speech and behaviour, yet it could also accommodate moments of the colloquial and the realistic. It encouraged a drama of range and versatility, with rapid changes of mood and great flexibility of tone. It is not surprising that in the 20th century the West has seen a return to the open stage and that recent plays of Brechtian theatre and the theatre of the absurd seem composed for open staging.

The proscenium stage. The third basic theatre form is that of the proscenium-arch or picture-frame stage, which reached its highest achievements in the late 19th century. Not until public theatres were roofed, the actors withdrawn into the scene, and the stage artificially illuminated were conditions ripe in Western theatre for a new development of spectacle and illusion. This development had a revolutionary effect upon the literary drama. In the 18th and 19th centuries, plays were shaped into a new structure of acts and scenes, with intermissions to permit scene changes. Only recently has the development of lighting techniques encouraged a return to a more flexible episodic drama. Of more importance, the actor increasingly withdrew into the created illusion of the play, and his character became part of it. In the mid-19th century, when it was possible to dim the house lights, the illusion could be made virtually complete. At its best, stage illusion could produce the delicate naturalism of a Chekhovian family scene, into which the spectator was drawn by understanding, sympathy, and recognition; at its worst, the magic of spectacle and the necessary projection of the speech and acting in the largest picture-frame theatres produced a crude drama of sensation in which literary values had no place.

Audience expectations. It may be that the primary influence upon the conception and creation of a play is that of the audience. An audience allows a play to have only the emotion and meaning it chooses, or else it defends itself either by protest or by a closed mind. From the time the spectator began paying for his playgoing, during

the Renaissance, the audience more and more entered into the choice of the drama's subjects and their treatment. This is not to say that the audience was given no consideration earlier; even in medieval plays there were popular non-biblical roles such as Noah's wife, or Mak the sheepthief among the three shepherds, and the antic devils of the Harrowing of Hell in the English mystery cycles. Nor, in later times, did a good playwright always give the audience only what it expected—Shakespeare's *King Lear* (c. 1605), for example, in the view of many the world's greatest play, had its popular elements of folktale, intrigue, disguise, madness, clowning, blood, and horror; but each was turned by the playwright to the advantage of his theme.

Any examination of the society an audience represents must illuminate not only the cultural role of its theatre but also the content, genre, and style of its plays. The exceptionally aristocratic composition of the English Restoration audience, for example, illuminates the social game its comedy represented, and the middle class composition of the subsequent Georgian audience sheds light on the moralistic elements of its "sentimental" comedy. Not unrelated is the study of received ideas in the theatre. The widespread knowledge of simple Freudian psychology has undoubtedly granted a contemporary playwright like Tennessee Williams (1911–83) the license to invoke it for character motivation; and Brecht increasingly informed his comedies with Marxist thinking on the assumption that the audiences he wrote for would appreciate his dramatized argument. Things go wrong when the intellectual or religious background of the audience does not permit a shared experience, as when Jean-Paul Sartre (1905–80) could not persuade a predominantly Christian audience with an existentialist explanation for the action of his plays, or when T.S. Eliot (1888–1965) failed to persuade an audience accustomed to the conventions of drawing-room comedy that *The Cocktail Party* (1949) was a possible setting for Christian martyrdom. Good drama persuades before it preaches, but it can only begin where the audience begins.

Special audiences. A great variety of drama has been written for special audiences. Plays have been written for children, largely in the 20th century, though Nativity plays have always been associated with children both as performers and as spectators. These plays tend to be fanciful in conception, broad in characterization, and moralistic in intention. Nevertheless, the most famous of children's plays, James Barrie's *Peter Pan* (1904), implied that the young are no fools and celebrated children in their own right. Barrie submerged his point subtly beneath the fantasy, and his play is still regularly performed, while Maurice Maeterlinck's *Blue Bird* (1908) has disappeared from the repertory because of its weighty moral tone.

In the wider field of adult drama, the social class of the audience often accounts for a play's form and style. Court or aristocratic drama is readily distinguished from that of the popular theatre. The veneration in which the Nō drama was held in Japan derived in large part from the feudal ceremony of its presentation, and its courtly elements ensured its survival for an upper class and intellectual elite. Although much of it derived from the Nō, the flourishing of the Kabuki at the end of the 17th century is related to the rise of a new merchant and middle class audience, which encouraged the development of less esoteric drama. The popular plays of the Elizabethan public theatres, with their broader, more romantic subjects liberally spiced with comedy, are similarly to be contrasted with those of the private theatres. The boys' companies of the private theatres of Elizabethan London played for a better paying and more sophisticated audience, which favoured the satirical or philosophical plays of Thomas Middleton (1570?–1627), John Marston (1576–1634), and George Chapman (1559?–1634). Similarly today, in all Western dramatic media—stage, film, radio, and television—popular and "commercial" forms run alongside more "cultural" and avant-garde forms, so that the drama, which in its origins brought people together, now divides them. Whether the esoteric influences the popular theatre, or vice versa, is not clear, and research remains to be done on whether

Children's
drama

this dichotomy is good or bad for dramatic literature or the people it is written for.

THE RANGE OF DRAMATIC FORMS AND STYLES

Dramatic literature has a remarkable facility in bringing together elements from other performing and nonperforming arts: design and mime, dance and music, poetry and narrative. It may be that the dramatic impulse itself, the desire to recreate a picture of life for others through impersonation, is at the root of all the arts. Certainly, the performing arts continually have need of dramatic literature to support them. A common way of describing an opera, for example, is to say that it is a play set to music. In Wagner the music is continuous; in Verdi the music is broken into songs; in Mozart the songs are separated by recitative, a mixture of speech and song; while operettas and musical comedy consist of speech that breaks into song from time to time. All forms of opera, however, essentially dramatize a plot, even if the plot must be simplified on the operatic stage. This is because, in opera, musical conventions dominate the dramatic conventions, and the spectator who finds that the music spoils the play, or who finds that the play spoils the music, is one who has not accepted the special conventions of opera. Music is drama's natural sister; proof may be seen in the early religious music-drama of the Dionysiac festivals of Greece and the *mystères* of 14th-century France, as well as in the remarkable development of opera in 17th-century Italy spreading to the rest of the world. The librettist who writes the text of an opera, however, must usually subserve the composer, unless he is able to embellish his play with popular lyrics, as John Gay did in *The Beggar's Opera* (1728), or to work in exceptionally close collaboration with the composer, as Brecht did with Kurt Weill for his *Die Dreigroschenoper* (*The Threepenny Opera*, 1928).

Dance, with its modern, sophisticated forms of ballet, has also been traditionally associated with dramatic representation and has similarly changed its purpose from religious to secular. In ballet, the music is usually central, and the performance is conceived visually and aurally; hence, the writer does not play a dominant role. The scenario is prepared for dance and mime by the choreographer. The contemporary Irish writer Samuel Beckett, trying to reduce his dramatic statement to the barest essentials, "composed" two mimes entitled *Act Without Words I* and *II* (1957 and 1966), but this is exceptional.

In motion pictures, the script writer has a more important but still not dominant role. He usually provides a loose outline of dialogue, business, and camera work on which the director, his cameramen, and the cutting editor build the finished product. The director is usually the final artistic authority and the central creative mind in the process, and words are usually subordinate to the dynamic visual imagery. (This subject is developed at length in the article MOTION PICTURES.)

The media of radio and television both depend upon words in their drama to an extent that is not characteristic of the motion picture. Though these mass media have been dominated by commercial interests and other economic factors, they also have developed dramatic forms from the special nature of their medium. The writer of a radio play must acknowledge that the listener cannot see the actors but hears them in conditions of great intimacy. A radio script that stresses the suggestive, imaginative, or poetic quality of words and permits a more than conventional freedom with time and place can produce a truly poetic drama, perhaps making unobtrusive use of earlier devices like the chorus, the narrator, and the soliloquy: the outstanding example of radio drama is *Under Milk Wood* (1953), by the Welsh poet Dylan Thomas.

A similar kind of dramatic writing is the so-called readers' theatre, in which actors read or recite without decor before an audience. (This is not to be confused with "closet drama," often a dramatic poem that assumes dialogue form; e.g., Milton's *Samson Agonistes*, 1671, written without the intention of stage performance.) The essential discipline of the circuit of communication with an audience is what distinguishes drama as a genre, however many forms it has taken in its long history. (J.L.S.)

Comedy

The classic conception of comedy, which began with Aristotle in ancient Greece of the 4th century BC and persists through the present, holds that it is primarily concerned with man as a social being, rather than as a private person, and that its function is frankly corrective. The comic artist's purpose is to hold a mirror up to society to reflect its follies and vices, in the hope that they will, as a result, be mended. The 20th-century French philosopher Henri Bergson shared this view of the corrective purpose of laughter: specifically, he felt, laughter is intended to bring the comic character back into conformity with his society, whose logic and conventions he abandons when "he slackens in the attention that is due to life." Here comedy is considered primarily as a literary genre, but also is considered for its manifestations in the other arts. The wellsprings of comedy are dealt with in the article HUMOUR AND WIT. The comic impulse in the visual arts is discussed in CARICATURE, CARTOON, AND COMIC STRIP.

ORIGINS AND DEFINITIONS

The word comedy seems to be connected by derivation with the Greek verb meaning "to revel," and comedy arose out of the revels associated with the rites of Dionysus, a god of vegetation. The origins of comedy are thus bound up with vegetation ritual. Aristotle, in his *Poetics*, states that comedy originated in phallic songs and that, like tragedy, it began in improvisation. Though tragedy evolved by stages that can be traced, the progress of comedy passed unnoticed because it was not taken seriously. When tragedy and comedy arose, poets wrote one or the other, according to their natural bent. Those of the graver sort, who might previously have been inclined to celebrate the actions of the great in epic poetry, turned to tragedy; poets of a lower type, who had set forth the doings of the ignoble in invectives, turned to comedy. The distinction is basic to the Aristotelian differentiation between tragedy and comedy: tragedy imitates men who are better than the average, and comedy men who are worse.

For centuries, efforts at defining comedy were to be along the lines set down by Aristotle: the view that tragedy deals with personages of high estate, and comedy deals with lowly types; that tragedy treats of matters of great public import, while comedy is concerned with the private affairs of mundane life; and that the characters and events of tragedy are historic and so, in some sense, true, while the humbler materials of comedy are but feigned. Implicit, too, in Aristotle is the distinction in styles deemed appropriate to the treatment of tragic and comic story. As long as there was at least a theoretical separation of comic and tragic styles, either genre could, on occasion, appropriate the stylistic manner of the other to a striking effect, which was never possible after the crossing of stylistic lines became commonplace. The ancient Roman poet Horace, who wrote on such stylistic differences, noted the special effects that can be achieved when comedy lifts its voice in pseudotragic rant and when tragedy adopts the prosaic but affecting language of comedy. Consciously combined, the mixture of styles produces the burlesque, in which the grand manner (epic or tragic) is applied to a trivial subject, or the serious subject is subjected to a vulgar treatment, to ludicrous effect. The English novelist Henry Fielding, in the preface to *Joseph Andrews* (1742), was careful to distinguish between the comic and the burlesque; the latter centres on the monstrous and unnatural and gives pleasure through the surprising absurdity it exhibits in appropriating the manners of the highest to the lowest, or vice versa. Comedy, on the other hand, confines itself to the imitation of nature, and, according to Fielding, the comic artist is not to be excused for deviating from it. His subject is the ridiculous, not the monstrous, as with the writer of burlesque; and the nature he is to imitate is human nature, as viewed in the ordinary scenes of civilized society.

The human contradiction. In dealing with man as a social being, all great comic artists have known that they are in the presence of a contradiction: that behind the social being lurks an animal being, whose behaviour of-

The corrective purpose of laughter

Distinction between the comic and the burlesque

ten accords very ill with the canons dictated by society. Comedy, from its ritual beginnings, has celebrated creative energy. The primitive revels out of which comedy arose frankly acknowledged man's animal nature; the animal masquerades and the phallic processions are the obvious witnesses to it. Comedy testifies to man's physical vitality, his delight in life, his will to go on living. Comedy is at its merriest, its most festive, when this rhythm of life can be affirmed within the civilized context of human society. In the absence of this sort of harmony between creatural instincts and the dictates of civilization, sundry strains and discontents arise, all baring witness to the contradictory nature of man, which in the comic view is a radical dualism; his efforts to follow the way of rational sobriety are forever being interrupted by the infirmities of the flesh. The duality that tragedy views as a fatal contradiction in the nature of things comedy views as one more instance of the incongruous reality that every man must live with as best he can. "Wherever there is life, there is contradiction," says Søren Kierkegaard, the 19th-century Danish Existentialist, in the *Concluding Unscientific Postscript* (1846), "and wherever there is contradiction, the comical is present." He went on to say that the tragic and the comic are both based on contradiction; but "the tragic is the suffering contradiction, comical, painless contradiction." Comedy makes the contradiction manifest along with a way out, which is why the contradiction is painless. Tragedy, on the other hand, despairs of a way out of the contradiction.

The incongruous is "the essence of the laughable," said the English essayist William Hazlitt, who also declared, in his essay "On Wit and Humour" in *English Comic Writers* (1819), that "Man is the only animal that laughs and weeps; for he is the only animal that is struck with the difference between what things are, and what they ought to be."

Comedy, satire, and romance. Comedy's dualistic view of man as an incongruous mixture of bodily instinct and rational intellect is an essentially ironic view—implying the capacity to see things in a double aspect. The comic drama takes on the features of satire as it fixes on professions of virtue and the practices that contradict them. Satire assumes standards against which professions and practices are judged. To the extent that the professions prove hollow and the practices vicious, the ironic perception darkens and deepens. The element of the incongruous points in the direction of the grotesque, which implies an admixture of elements that do not match. The ironic gaze eventually penetrates to a vision of the grotesque quality of experience, marked by the discontinuity of word and deed and the total lack of coherence between appearance and reality. This suggests one of the extreme limits of comedy, the satiric extreme, in which the sense of the discrepancy between things as they are and things as they might be or ought to be has reached to the borders of tragedy. For the tragic apprehension, as Kierkegaard states, despairs of a way out of the contradictions that life presents.

As satire may be said to govern the movement of comedy in one direction, romance governs its movement in the other. Satiric comedy dramatizes the discrepancy between the ideal and the reality and condemns the pretensions that would mask reality's hollowness and viciousness. Romantic comedy also regularly presents the conflict between the ideal shape of things as hero or heroine could wish them to be and the hard realities with which they are confronted, but typically it ends by invoking the ideal, despite whatever difficulties reality has put in its way. This is never managed without a good deal of contrivance, and the plot of the typical romantic comedy is a medley of clever scheming, calculated coincidence, and wondrous discovery, all of which contribute ultimately to making the events answer precisely to the hero's or heroine's wishes. Plotting of this sort has had a long stage tradition and not exclusively in comedy. It is first encountered in the tragicomedies of the ancient Greek dramatist Euripides (e.g., *Alcestris*, *Iphigenia in Tauris*, *Ion*, *Helen*). Shakespeare explored the full range of dramatic possibilities of the romantic mode of comedy. The means by which the happy ending is accomplished in romantic comedy—the

document or the bodily mark that establishes identities to the satisfaction of all the characters of goodwill—are part of the stock-in-trade of all comic dramatists, even such 20th-century playwrights as Jean Anouilh (in *Le Voyageur sans bagage*) and T.S. Eliot (in *The Confidential Clerk*).

There is nothing necessarily inconsistent in the use of a calculatedly artificial dramatic design to convey a serious dramatic statement. The contrived artifice of Shakespeare's mature comic plots is the perfect foil against which the reality of the characters' feelings and attitudes assumes the greater naturalness. The strange coincidences, remarkable discoveries, and wonderful reunions are unimportant compared with the emotions of relief and awe that they inspire. Their function, as Shakespeare uses them, is precisely to give rise to such emotions, and the emotions, thanks to the plangent poetry in which they are expressed, end by transcending the circumstances that occasioned them. But when such artifices are employed simply for the purpose of eliminating the obstacles to a happy ending—as is the case in the sentimental comedy of the 18th and early 19th centuries—then they stand forth as imaginatively impoverished dramatic clichés. The dramatists of sentimental comedy were committed to writing exemplary plays, wherein virtue would be rewarded and vice frustrated. If hero and heroine were to be rescued from the distresses that had encompassed them, any measures were apparently acceptable; the important thing was that the play's action should reach an edifying end. It is but a short step from comedy of this sort to the melodrama that flourished in the 19th-century theatre. The distresses that the hero and heroine suffer are, in melodrama, raised to a more than comic urgency, but the means of deliverance have the familiar comic stamp: the secret at last made known, the long-lost child identified, the hard heart made suddenly capable of pity. Melodrama is a form of fantasy that proceeds according to its own childish and somewhat egoistic logic: hero and heroine are pure, anyone who opposes them is a villain, and the purity that has exposed them to risks must ensure their eventual safety and happiness. What melodrama is to tragedy farce is to comedy, and the element of fantasy is equally prominent in farce and in melodrama. If melodrama provides a fantasy in which the protagonist suffers for his virtues but is eventually rewarded for them, farce provides a fantasy in which the protagonist sets about satisfying his most roguish or wanton, mischievous or destructive, impulses and manages to do so with impunity.

THEORIES

The treatise that Aristotle is presumed to have written on comedy is lost. There is, however, a fragmentary treatise on comedy that bears an obvious relation to Aristotle's treatise on tragedy, *Poetics*, and is generally taken to be either a version of a lost Aristotelian original or an expression of the philosophical tradition to which he belonged. This is the *Tractatus Coislinianus*, preserved in a 10th-century manuscript in the De Coislin Collection in Paris. The *Tractatus* divides the substance of comedy into the same six elements that are discussed in regard to tragedy in the *Poetics*: plot, character, thought, diction, melody, and spectacle. The characters of comedy, according to the *Tractatus*, are of three kinds: the impostors, the self-deprecators, and the buffoons. The Aristotelian tradition from which the *Tractatus* derives probably provided a fourth, the churl, or boor. The list of comic characters in the *Tractatus* is closely related to a passage in Aristotle's *Nicomachean Ethics*, in which the boaster (the person who says more than the truth) is compared with the mock-modest man (the person who says less), and the buffoon (who has too much wit) is contrasted with the boor (who has too little).

Comedy as a rite. The *Tractatus* was not printed until 1839, and its influence on comic theory is thus of relatively modern date. It is frequently cited in the studies that attempt to combine literary criticism and anthropology, in the manner in which Sir James George Frazer combined studies of primitive religion and culture in *The Golden Bough* (1890–1915). In such works, comedy and tragedy alike are traced to a prehistoric death-and-resurrection

The ironic
view
of man

19th-
century
melodrama

ceremonial, a seasonal pantomime in which the old year, in the guise of an aged king (or hero or god), is killed, and the new spirit of fertility, the resurrection or initiation of the young king, is brought in. This rite typically featured a ritual combat, or agon, between the representatives of the old and the new seasons, a feast in which the sacrificial body of the slain king was devoured, a marriage between the victorious new king and his chosen bride, and a final triumphal procession in celebration of the reincarnation or resurrection of the slain god. Implicit in the whole ceremony is the ancient rite of purging the tribe through the expulsion of a scapegoat, who carries away the accumulated sins of the past year. Frazer, speaking of scapegoats in *The Golden Bough*, noted that this expulsion of devils was commonly preceded or followed by a period of general license, an abandonment of the ordinary restraints of society during which all offenses except the gravest go unpunished. This quality of Saturnalia is characteristic of comedy from ancient Greece through medieval Europe.

The seasonal rites that celebrate the yearly cycle of birth, death, and rebirth are seen by the contemporary Canadian critic Northrop Frye as the basis for the generic plots of comedy, romance, tragedy, and irony and satire. The four prefigure the fate of a hero and the society he brings into being. In comedy (representing the season of spring), the hero appears in a society controlled by obstructing characters and succeeds in wresting it from their grasp. The movement of comedy of this sort typically replaces falsehood with truth, illusion with reality. The hero, having come into possession of his new society, sets forth upon adventures, and these are the province of romance (summer). Tragedy (autumn) commemorates the hero's passion and death. Irony and satire (winter) depict a world from which the hero has disappeared, a vision of "unidealized existence." With spring, the hero is born anew.

The moral force of comedy. The characters of comedy specified in the *Tractatus* arrange themselves in a familiar pattern: a clever hero is surrounded by fools of sundry varieties (impostors, buffoons, boors). The hero is something of a trickster; he dissimulates his own powers, while exploiting the weaknesses of those around him. The comic pattern is a persistent one; it appears not only in ancient Greek comedy but also in the farces of ancient Italy, in the commedia dell'arte that came into being in 16th-century Italy, and even in the routines involving a comedian and his straight man in the nightclub acts and the television variety shows of the present time. Implicit here is the tendency to make folly ridiculous, to laugh it out of countenance, which has always been a prominent feature of comedy.

Renaissance critics, elaborating on the brief and cryptic account of comedy in Aristotle's *Poetics*, stressed the derisive force of comedy as an adjunct to morality. The Italian scholar Gian Giorgio Trissino's account of comedy in his *Poetica*, apparently written in the 1530s, is typical: as tragedy teaches by means of pity and fear, comedy teaches by deriding things that are vile. Attention is directed here, as in other critical treatises of this kind, to the source of laughter. According to Trissino, laughter is aroused by objects that are in some way ugly and especially by that from which better qualities were hoped. His statement suggests the relation of the comic to the incongruous. Trissino was as aware as the French poet Charles Baudelaire was three centuries later that laughter betokens the fallen nature of man (Baudelaire would term it man's Satanic nature). Man laughs, says Trissino (echoing Plato's dialogue *Philebus*), because he is envious and malicious and never delights in the good of others except when he hopes for some good from it for himself.

The most important English Renaissance statement concerning comedy is that of Sir Philip Sidney in *The Defence of Poesie* (1595):

comedy is an imitation of the common errors of our life, which [the comic dramatist] representeth in the most ridiculous and scornful sort that may be, so as it is impossible that any beholder can be content to be such a one.

Like Trissino, Sidney notes that, while laughter comes from delight, not all objects of delight cause laughter,

and he demonstrates the distinction as Trissino had done: "we are ravished with delight to see a fair woman, and yet are far from being moved to laughter. We laugh at deformed creatures, wherein certainly we cannot delight." The element of the incongruous is prominent in Sidney's account of scornful laughter. He cites the image of the hero of Greek legend Heracles, with his great beard and furious countenance, in woman's attire, spinning at the command of his beloved queen, Omphale, and declares that this arouses both delight and laughter.

Comedy and character. Another English poet, John Dryden, in *Of Dramatick Poesie, an Essay* (1668), makes the same point in describing the kind of laughter produced by the ancient Greek comedy *The Clouds*, by Aristophanes. In it, the character of Socrates is made ridiculous by acting very unlike the true Socrates; that is, by appearing childish and absurd rather than with the gravity of the true Socrates. Dryden was concerned with analyzing the laughable quality of comedy and with demonstrating the different forms it has taken in different periods of dramatic history. Aristophanic comedy sought its laughable quality not so much in the imitation of a man as in the representation of "some odd conceit which had commonly somewhat of unnatural or obscene in it." In the so-called New Comedy, introduced by Menander late in the 4th century BC, writers sought to express the ethos, or character, as in their tragedies they expressed the pathos, or suffering, of mankind. This distinction goes back to Aristotle, who, in the *Rhetoric*, distinguished between ethos, a man's natural bent, disposition, or moral character, and pathos, emotion displayed in a given situation. And the Latin rhetorician Quintilian, in the 1st century AD, noted that ethos is akin to comedy and pathos to tragedy. The distinction is important to Renaissance and Neoclassical assumptions concerning the respective subject of comic and tragic representation. In terms of emotion, ethos is viewed as a permanent condition characteristic of the average man and relatively mild in its nature; pathos, on the other hand, is a temporary emotional state, often violent. Comedy thus expresses the characters of men in the ordinary circumstances of everyday life; tragedy expresses the sufferings of a particular man in extraordinary periods of intense emotion.

In dealing with men engaged in normal affairs, the comic dramatists tended to depict the individual in terms of some single but overriding personal trait or habit. They adopted a method based on the physiological concept of the four humours, or bodily fluids (blood, phlegm, cholera, melancholy), and the belief that an equal proportion of these constituted health, while an excess or deficiency of any one of them brought disease. Since the humours governed temperament, an irregular distribution of them was considered to result not only in bodily sickness but also in derangements of personality and behaviour, as well. The resultant comedy of humours is distinctly English, as Dryden notes, and particularly identified with the comedies of Ben Jonson.

The role of wit. Humour is native to man. Folly need only be observed and imitated by the comic dramatist to give rise to laughter. Observers as early as Quintilian, however, have pointed out that, though folly is laughable in itself, such jests may be improved if the writer adds something of his own; namely, wit. A form of repartee, wit implies both a mental agility and a linguistic tact that is very much a product of conscious art. Quintilian describes wit at some length in his *Institutio oratoria*; it partakes of urbanity, a certain tincture of learning, charm, saltiness, or sharpness, and polish and elegance. In the preface (1671) to *An Evening's Love*, Dryden distinguishes between the comic talents of Ben Jonson, on the one hand, and of Shakespeare and his contemporary John Fletcher, on the other, by virtue of their excelling, respectively, in humour and wit. Jonson's talent lay in his ability "to make men appear pleasantly ridiculous on the stage"; while Shakespeare and Fletcher excelled in wit, or "the sharpness of conceit," as seen in their repartee. The distinction is noted as well in *Of Dramatick Poesie, an Essay*, where a comparison is made between the character of Morose in Jonson's play *Epitaph*, who is characterized by his

The role of the scapegoats

Sir Philip Sidney's definition of comedy

English comedy of humours

humour (namely, his inability to abide any noise but the sound of his own voice), and Shakespeare's Falstaff, who, according to Dryden, represents a miscellany of humours and is singular in saying things that are unexpected by the audience.

The distinctions that Hazlitt arrives at, then, in his essay "On Wit and Humour" are very much in the classic tradition of comic criticism:

Humour is the describing the ludicrous as it is in itself; wit is the exposing it, by comparing or contrasting it with something else. Humour is, as it were, the growth of nature and accident; wit is the product of art and fancy.

The distinctions persist into the most sophisticated treatments of the subject. Sigmund Freud, for example, in *Wit and its Relation to the Unconscious* (1905), said that wit is made, but humour is found. Laughter, according to Freud, is aroused at actions that appear immoderate and inappropriate, at excessive expenditures of energy: it expresses a pleasurable sense of the superiority felt on such occasions.

Baudelaire on the grotesque. The view that laughter comes from superiority is referred to as a commonplace by Baudelaire, who states it in his essay "On the Essence of Laughter" (1855). Laughter, says Baudelaire, is a consequence of man's notion of his own superiority. It is a token both of an infinite misery, in relation to the absolute being of whom man has an inkling, and of infinite grandeur, in relation to the beasts, and results from the perpetual collision of these two infinities. The crucial part of Baudelaire's essay, however, turns on his distinction between the comic and the grotesque. The comic, he says, is an imitation mixed with a certain creative faculty; the grotesque is a creation mixed with a certain imitative faculty—imitative of elements found in nature. Each gives rise to laughter expressive of an idea of superiority—in the comic, the superiority of man over man, and, in the grotesque, the superiority of man over nature. The laughter caused by the grotesque has about it something more profound and primitive, something much closer to the innocent life, than has the laughter caused by the comic in man's behaviour. In France, the great master of the grotesque was the 16th-century author François Rabelais, while some of the plays of Molière, in the next century, best expressed the comic.

Bergson's and Meredith's theories. The French philosopher Henri Bergson (1859–1941) analyzed the dialectic of comedy in his essay "Laughter," which deals directly with the spirit of contradiction that is basic both to comedy and to life. Bergson's central concern is with the opposition of the mechanical and the living; stated in its most general terms, his thesis holds that the comic consists of something mechanical encrusted on the living. Bergson traces the implications of this view in the sundry elements of comedy: situations, language, characters. Comedy expresses a lack of adaptability to society; any individual is comic who goes his own way without troubling to get into touch with his fellow beings. The purpose of laughter is to wake him from his dream. Three conditions are essential for the comic: the character must be unsociable, for that is enough to make him ludicrous; the spectator must be insensible to the character's condition, for laughter is incompatible with emotion; and the character must act automatically (Bergson cites the systematic absentmindedness of Don Quixote). The essential difference between comedy and tragedy, says Bergson, invoking a distinction that goes back to that maintained between ethos and pathos, is that tragedy is concerned with individuals and comedy with classes. And the reason that comedy deals with the general is bound up with the corrective aim of laughter: the correction must reach as great a number of persons as possible. To this end, comedy focusses on peculiarities that are not indissolubly bound up with the individuality of a single person.

It is the business of laughter to repress any tendency on the part of the individual to separate himself from society. The comic character would, if left to his own devices, break away from logic (and thus relieve himself from the strain of thinking); give over the effort to adapt and readapt himself to society (and thus slacken in the atten-

tion that is due to life); and abandon social convention (and thus relieve himself from the strain of living).

The essay "On the Idea of Comedy and the Uses of the Comic Spirit" (1877), by Bergson's English contemporary George Meredith, is a celebration of the civilizing power of the comic spirit. The mind, he affirms, directs the laughter of comedy, and civilization is founded in common sense, which equips one to hear the comic spirit when it laughs folly out of countenance and to participate in its fellowship.

Both Bergson's and Meredith's essays have been criticized for focussing so exclusively on comedy as a socially corrective force and for limiting the scope of laughter to its derisive power. The charge is more damaging to Meredith's essay than it is to Bergson's. Whatever the limitations of the latter, it nonetheless explores the implications of its own thesis with the utmost thoroughness, and the result is a rigorous analysis of comic causes and effects for which any student of the subject must be grateful. It is with farce that Bergson's remarks on comedy have the greatest connection and on which they seem chiefly to have been founded. It is no accident that most of his examples are drawn from Molière, in whose work the farcical element is strong, and from the farces of Bergson's own contemporary Eugène Labiche. The laughter of comedy is not always derisive, however, as some of Shakespeare's greatest comedies prove; and there are plays, such as Shakespeare's last ones, which are well within an established tradition of comedy but in which laughter hardly sounds at all. These suggest regions of comedy on which Bergson's analysis of the genre sheds hardly any light at all.

The comic as a failure of self-knowledge. Aristotle said that comedy deals with the ridiculous, and Plato, in the *Philebus*, defined the ridiculous as a failure of self-knowledge: such a failure is there shown to be laughable in private individuals (the personages of comedy) but terrible in persons who wield power (the personages of tragedy). In comedy, the failure is often mirrored in a character's efforts to live up to an ideal of self that may be perfectly worthy but the wrong ideal for him. Shakespearean comedy is rich in examples: the King of Navarre and his courtiers, who must be made to realize that nature meant them to be lovers, not academicians, in *Love's Labour's Lost*; Beatrice and Benedick, who must be made to know that nature meant them for each other, not for the single life, in *Much Ado About Nothing*; the Duke Orsino in *Twelfth Night*, who is brought to see that it is not Lady Olivia whom he loves but the disguised Viola, and Lady Olivia herself, who, when the right man comes along, decides that she will not dedicate herself to seven years of mourning for a dead brother, after all; and Angelo in *Measure for Measure*, whose image of himself collapses when his lust for Isabella makes it clear that he is not the ascetic type. The movement of all these plays follows a familiar comic pattern, wherein characters are brought from a condition of affected folly amounting to self-delusion to a plain recognition of who they are and what they want. For the five years or so after he wrote *Measure for Measure*, in 1604, Shakespeare seems to have addressed himself exclusively to tragedy, and each play in the sequence of masterpieces he produced during this period—*Othello*, *King Lear*, *Macbeth*, *Antony and Cleopatra*, and *Coriolanus*—turns in some measure on a failure of self-knowledge. This is notably so in the case of *Lear*, which is the tragedy of a man who (in the words of one of his daughters) "hath ever but slenderly known himself," and whose fault (as the Fool suggests) is to have grown old before he grew wise.

The plots of Shakespeare's last plays (*Pericles*, *Cymbeline*, *The Winter's Tale*, *The Tempest*) all contain a potential tragedy but one that is resolved by nontragic means. They contain, as well, an element of romance of the kind purveyed from Greek New Comedy through the plays of the ancient Roman comic dramatists Plautus and Terence. Children lost at birth are miraculously restored, years later, to their parents, thereby providing occasion for a recognition scene that functions as the denouement of the plot. Characters find themselves—they come to know themselves—in all manner of ways by the ends of these

Baudelaire's definition of laughter

The laughter of comedy

plays. Tragic errors have been made, tragic losses have been suffered, tragic passions—envy, jealousy, wrath—have seemed to rage unchecked, but the miracle that these plays celebrate lies in the discovery that the errors can be forgiven, the losses restored, and the passions mastered by man's godly spirit of reason. The near tragedies experienced by the characters result in the ultimate health and enlightenment of the soul. What is learned is of a profound simplicity: the need for patience under adversity, the need to repent of one's sins, the need to forgive the sins of others. In comedy of this high and sublime sort, patience, repentance, and forgiveness are opposed to the viciously circular pattern of crime, which begets vengeance, which begets more crime. Comedy of this sort deals in regeneration and rebirth. There is always about it something of the religious, as humankind is absolved of its guilt and reconciled one to another and to whatever powers that be.

Divine comedies in the West and East. The 4th-century Latin grammarian Donatus distinguished comedy from tragedy by the simplest terms: comedies begin in trouble and end in peace, while tragedies begin in calms and end in tempest. Such a differentiation of the two genres may be simplistic, but it provided sufficient grounds for Dante to call his great poem *La Commedia (The Comedy)*; later called *The Divine Comedy*, since, as he says in his dedicatory letter, it begins amid the horrors of hell but ends amid the pleasures of heaven. This suggests the movement of Shakespeare's last plays, which begin amid the distresses of the world and end in a supernal peace. Comedy conceived in this sublime and serene mode is rare but recurrent in the history of the theatre. The Spanish dramatist Calderón's *Vida es sueño* (1635; "Life Is a Dream") is an example; so, on the operatic stage, is Mozart's *Magic Flute* (1791), in spirit and form so like Shakespeare's *Tempest*, to which it has often been compared. In later drama, Henrik Ibsen's *Little Eyolf* (1894) and August Strindberg's *To Damascus* (1898–1904)—both of which are among the late works of these Scandinavian dramatists—have affinities with this type, and this is the comic mode in which T.S. Eliot's last play, *The Elder Statesman* (1958), is conceived. It may represent the most universal mode of comedy. The American philosopher Susanne K. Langer writes:

In Asia the designation "Divine Comedy" would fit numberless plays; especially in India triumphant gods, divine lovers united after various trials [as in the perennially popular romance of Rama and Sita], are the favourite themes of a theater that knows no "tragic rhythm." The classical Sanskrit drama was heroic comedy—high poetry, noble action, themes almost always taken from the myths—a serious, religiously conceived drama, yet in the "comic" pattern, which is not a complete organic development reaching a foregone, inevitable conclusion, but is episodic, restoring a lost balance, and implying a new future. The reason for this consistently "comic" image of life in India is obvious enough: both Hindu and Buddhist regard life as an episode in the much longer career of the soul which has to accomplish many incarnations before it reaches its goal, nirvana. Its struggles in the world do not exhaust it; in fact they are scarcely worth recording except in entertainment theater, "comedy" in our sense—satire, farce, and dialogue. The characters whose fortunes are seriously interesting are the eternal gods; and for them there is no death, no limit of potentialities, hence no fate to be fulfilled. There is only the balanced rhythm of sentience and emotion, upholding itself amid the changes of material nature. (From *Feeling and Form*; Charles Scribner's Sons, 1953.)

KINDS OF COMEDY IN DIVERSE HISTORICAL PERIODS

Old and New Comedy in ancient Greece. The 11 surviving plays of Aristophanes represent the earliest extant body of comic drama; what is known of Greek Old Comedy is derived from these plays, the earliest of which, *The Acharnians*, was produced in 425 BC. Aristophanic comedy has a distinct formal design but displays very little plot in any conventional sense. Rather, it presents a series of episodes aimed at illustrating, in humorous and often bawdy detail, the implications of a deadly serious political issue: it is a blend of invective, buffoonery, and song and dance. Old Comedy often used derision and scurrility, and this may have proved its undoing; though praised by all, the freedom it enjoyed degenerated into license and violence and had to be checked by law.

In New Comedy, which began to prevail around 336 BC, the Aristophanic depiction of public personages and events was replaced by a representation of the private affairs (usually amorous) of imaginary men and women. New Comedy is known only from the fragments that have survived of the plays of Menander (c. 342–c. 292 BC) and from plays written in imitation of the form by the Romans Plautus (c. 254–184 BC) and Terence (195 or 185–159 BC). A number of the stock comic characters survived from Old Comedy into New: an old man, a young man, an old woman, a young woman, a learned doctor or pedant, a cook, a parasite, a swaggering soldier, a comic slave. New Comedy, on the other hand, exhibits a degree of plot articulation never achieved in the Old. The action of New Comedy is usually about plotting; a clever servant, for example, devises ingenious intrigues in order that his young master may win the girl of his choice. There is satire in New Comedy: on a miser who loses his gold from being overcareful of it (the *Aulularia* of Plautus); on a father who tries so hard to win the girl from his son that he falls into a trap set for him by his wife (Plautus' *Casina*); and on an over stern father whose son turns out worse than the product of an indulgent parent (in the *Adelphi* of Terence). But the satiric quality of these plays is bland by comparison with the trenchant ridicule of Old Comedy. The emphasis in New Comic plotting is on the conduct of a love intrigue; the love element per se is often of the slightest, the girl whom the hero wishes to possess sometimes being no more than an offstage presence or, if onstage, a mute.

New Comedy provided the model for European comedy through the 18th century. During the Renaissance, the plays of Plautus and, especially, of Terence were studied for the moral instruction that young men could find in them: lessons on the need to avoid the snares of harlots and the company of braggarts, to govern the deceitful trickery of servants, to behave in a seemly and modest fashion to parents. Classical comedy was brought up to date in the plays of the "Christian Terence," imitations by schoolmasters of the comedies of the Roman dramatist. They added a contemporary flavour to the life portrayed and displayed a somewhat less indulgent attitude to youthful indiscretions than did the Roman comedy. New Comedy provided the basic conventions of plot and characterization for the *commedia erudita*—comedy performed from written texts—of 16th-century Italy, as in the plays of Machiavelli and Ariosto. Similarly, the stock characters that persisted from Old Comedy into New were taken over into the improvisational *commedia dell'arte*, becoming such standard masked characters as Pantalone, the Dottore, the vainglorious Capitano, the young lovers, and the servants, or *zanni*.

Rise of realistic comedy in 17th-century England. The early part of the 17th century in England saw the rise of a realistic mode of comedy based on a satiric observation of contemporary manners and mores. It was masterminded by Ben Jonson, and its purpose was didactic. Comedy, said Jonson in *Every Man Out of his Humour* (1599), quoting the definition that during the Renaissance was attributed to Cicero, is an imitation of life, a glass of custom, an image of truth. Comedy holds the mirror up to nature and reflects things as they are, to the end that society may recognize the extent of its shortcomings and the folly of its ways and set about its improvement. Jonson's greatest plays—*Volpone* (1606), *Epicure* (1609), *The Alchemist* (1610), *Bartholomew Fair* (1614)—offer a richly detailed contemporary account of the follies and vices that are always with us. The setting (apart from *Volpone*) is Jonson's own London, and the characters are the ingenious or the devious or the grotesque products of the human wish to get ahead in the world. The conduct of a Jonsonian comic plot is in the hands of a clever manipulator who is out to make reality conform to his own desires. Sometimes he succeeds, as in the case of the clever young gentleman who gains his uncle's inheritance in *Epicure* or the one who gains the rich Puritan widow for his wife in *Bartholomew Fair*. In *Volpone* and *The Alchemist*, the schemes eventually fail, but this is the fault of the manipulators, who will never stop when they are ahead, and not at all due to any

Distinction between comedy and tragedy

Development of stock characters

The rise of a realistic comedy

insight on the part of the victims. The victims are almost embarrassingly eager to be victimized. Each has his ruling passion—his humour—and it serves to set him more or less mechanically in the path that he will undeviatingly pursue, to his own discomfiture.

English comedy of the later 17th century is cast in the Jonsonian mold. Restoration comedy is always concerned with the same subject—the game of love—but the subject is treated as a critique of fashionable society. Its aim is distinctly satiric, and it is set forth in plots of Jonsonian complexity, where the principal intriguer is the rakish hero, bent on satisfying his sexual needs, outside the bonds of marriage, if possible. In the greatest of these comedies—Sir George Etherege's *Man of Mode* (1676), for example, or William Wycherley's *Country-Wife* (1675) or William Congreve's *Way of the World* (1700)—the premium is on the energy and the grace with which the game is played, and the highest dramatic approval is reserved for those who take the game seriously enough to play it with style but who have the good sense to know when it is played out. The satiric import of Restoration comedy resides in the dramatist's awareness of a familiar incongruity: that between the image of man in his primitive nature and the image of man amid the artificial restraints that society would impose upon him. The satirist in these plays is chiefly concerned with detailing the artful dodges that ladies and gentlemen employ to satisfy nature and to remain within the pale of social decorum. Inevitably, then, hypocrisy is the chief satiric target. The animal nature of man is taken for granted, and so is his social responsibility to keep up appearances; some hypocrisy must follow, and, within limits, society will wink at indiscretions so long as they are discreetly managed. The paradox is typical of those in which the Restoration comic dramatists delight: and the strongly rational and unidealistic ethos of this comedy has its affinities with the naturalistic and skeptical cast of late-17th-century philosophical thought.

Sentimental comedy of the 17th and 18th centuries. The Restoration comic style collapsed around the end of the 17th century, when the satiric vision gave place to a sentimental one. Jeremy Collier's *Short view of the Profaneness and Immorality of the English Stage*, published in 1698, signalled the public opposition to the real or fancied improprieties of plays staged during the previous three decades. "The business of plays is to recommend Vertue, and discountenance Vice"; so runs the opening sentence of Collier's attack. No Restoration comic dramatist ever conceived of his function in quite these terms. "It is the business of a comic poet to paint the vices and follies of humankind," Congreve had written a few years earlier (in the dedication to *The Double-Dealer*). Though Congreve may be assumed to imply—in accordance with the time-honoured theory concerning the didactic end of comedy—that the comic dramatist paints the vices and follies of humankind for the purpose of correcting them through ridicule, he is, nonetheless, silent on this point. Collier's assumption that all plays must recommend virtue and discountenance vice has the effect of imposing on comedy the same sort of moral levy that critics such as Thomas Rymer were imposing on tragedy in their demand that it satisfy poetic justice.

At the beginning of the 18th century, there was a blending of the tragic and comic genres that, in one form or another, had been attempted throughout the preceding century. The vogue of tragicomedy may be said to have been launched in England with the publication of John Fletcher's *Faithfull Shepheardesse* (c. 1608), an imitation of the *Pastor fido*, by the Italian poet Battista Guarini. In his *Compendium of Tragicomic Poetry* (1601), Guarini had argued the distinct nature of the genre, maintaining it to be a third poetic kind, different from either the comic or the tragic. Tragicomedy, he wrote, takes from tragedy its great persons but not its great action, its movement of the feelings but not its disturbance of them, its pleasure but not its sadness, its danger but not its death; from comedy it takes laughter that is not excessive, modest amusement, feigned difficulty, and happy reversal. Fletcher adapted this statement in the address "To the Reader" that prefaces *The Faithfull Shepheardesse*.

The form quickly established itself on the English stage, and, through the force of such examples as Beaumont and Fletcher's *Phylaster* (1610) and *A King and No King* (1611) and a long sequence of Fletcher's unaided tragicomedies, it prevailed during the 20 years before the closing of the theatres in 1642. The taste for tragicomedy continued unabated at the Restoration, and its influence was so pervasive that during the closing decades of the century the form began to be seen in plays that were not, at least by authorial designation, tragicomedies. Its effect on tragedy can be seen not only in the tendency, always present on the English stage, to mix scenes of mirth with more solemn matters but also in the practice of providing tragedy with a double ending (a fortunate one for the virtuous, an unfortunate one for the vicious), as in Dryden's *Aureng-Zebe* (1675) or Congreve's *Mourning Bride* (1697). The general lines separating the tragic and comic genres began to break down, and that which is high, serious, and capable of arousing pathos could exist in the same play with what is low, ridiculous, and capable of arousing derision. The next step in the process came when Sir Richard Steele, bent on reforming comedy for didactic purposes, produced *The Conscious Lovers* (1722) and provided the English stage with an occasion when the audience at a comedy could derive its chief pleasure not from laughing but from weeping. It wept in the delight of seeing virtue rewarded and young love come to flower after parental opposition had been overcome. Comedy of the sort inaugurated by *The Conscious Lovers* continued to represent the affairs of private life, as comedy had always done, but with a seriousness hitherto unknown; and the traditionally low personages of comedy now had a capacity for feeling that bestowed on them a dignity previously reserved for the personages of tragedy.

This trend in comedy was part of a wave of egalitarianism that swept through 18th-century political and social thought. It was matched by a corresponding trend in tragedy, which increasingly selected its subjects from the affairs of private men and women in ordinary life, rather than from the doings of the great. The German dramatist Gotthold Lessing wrote that the misfortunes of those whose circumstances most resemble those of the audience most naturally penetrate most deeply into its heart, and his own *Minna von Barnhelm* (1767) is an example of the new serious comedy. The capacity to feel, to sympathize with, and to be affected by the plight of a fellow human being without regard for his rank in the world's esteem became the measure of one's humanity. It was a bond that united the fraternity of mankind in an aesthetic revolution that preceded the political revolutions of the 18th century. In literature, this had the effect of hastening the movement toward a more realistic representation of reality, whereby the familiar events of common life are treated "seriously and problematically" (in the phrase of the critic Erich Auerbach, who traced the process in his book *Mimesis* [1946]). The results may be seen in novels such as Samuel Richardson's *Pamela* and *Clarissa* and in middle-class tragedies such as George Lillo's *The London Merchant* (1731) in England; in the *comédie larmoyante* ("tearful comedy") in France; in Carlo Goldoni's efforts to reform the *commedia dell'arte* and replace it with a more naturalistic comedy in the Italian theatre; and in the English sentimental comedy, exemplified in its full-blown state by plays such as Hugh Kelly's *False Delicacy* (1768) and Richard Cumberland's *West Indian* (1771). Concerning the sentimental comedy it must be noted that it is only in the matter of appropriating for the bourgeoisie a seriousness of tone and a dignity of representational style previously considered the exclusive property of the nobility that the form can be said to stand in any significant relationship to the development of a more realistic mimetic mode than the traditional tragic and comic ones. The plots of sentimental comedy are as contrived as anything in Plautus and Terence (which with their fondness for foundling heroes who turn out to be long-lost sons of rich merchants, they often resemble); and with their delicate feelings and genteel moral atmosphere, comedies of this sort seem as affected in matters of sentiment as Restoration comedy seems in matters of wit.

The trend to egalitarianism

The vogue of tragicomedy

Goldsmith's views on sentimental comedy

Oliver Goldsmith, in his "A Comparison Between Laughing and Sentimental Comedy" (1773), noted the extent to which the comedy in the England of his day had departed from its traditional purpose, the excitation of laughter by exhibiting the follies of the lower part of mankind. He questioned whether an exhibition of its follies would not be preferable to a detail of its calamities. In sentimental comedy, Goldsmith continued, the virtues of private life were exhibited, rather than the vices exposed; and the distresses rather than the faults of mankind generated interest in the piece. Characters in these plays were almost always good; if they had faults, the spectator was expected not only to pardon but to applaud them, in consideration of the goodness of their hearts. Thus, according to Goldsmith, folly was commended instead of being ridiculed. Goldsmith concluded by labelling sentimental comedy a "species of bastard tragedy," "a kind of *mulish* production": a designation that ironically brings to mind Guarini's comparison of tragicomedy in its uniqueness (a product of comedy and tragedy but different from either) to the mule (the offspring of the horse and the ass but itself neither one nor the other). The production of Goldsmith's *She Stoops to Conquer* (1773) and of Richard Brinsley Sheridan's *Rivals* (1775) and *The School for Scandal* (1777) briefly reintroduced comic gaiety to the English stage; by the end of the decade, Sheridan's dramatic burlesque, *The Critic* (first performed 1779), had appeared, with its parody of contemporary dramatic fashions, the sentimental included. But this virtually concluded Sheridan's career as a dramatist; Goldsmith had died in 1774; and the sentimental play was to continue to govern the English comic stage for over a century to come.

The comic outside the theatre. The great comic voices of the 18th century in England were not those in the theatre. No dramatic satire of the period can exhibit anything comparable to the furious ridicule of man's triviality and viciousness that Jonathan Swift provided in *Gulliver's Travels* (1726). His *Modest Proposal* (1729) is a masterpiece of comic incongruity, with its suave blend of rational deliberation and savage conclusion. The comic artistry of Alexander Pope is equally impressive. Pope expressed his genius in the invective of his satiric portraits and in the range of moral and imaginative vision that was capable, at one end of his poetic scale, of conducting that most elegant of drawing-room epics, *The Rape of the Lock* (1712–14), to its sublimely inane conclusion and, at the other, of invoking from the scene that closes *The Dunciad* (1728) an apocalyptic judgment telling what will happen when the vulgarizers of the word have carried the day.

When the voice of comedy did sound on the 18th-century English stage with anything approaching its full critical and satiric resonance, the officials soon silenced it. John Gay's *Beggar's Opera* (1728) combined hilarity with a satiric fierceness worthy of Swift (who may have suggested the original idea for it). The officials tolerated its spectacularly successful run, but no license from the lord chamberlain could be secured for Gay's sequel, *Polly*, which was not staged until 1777. The Licensing Act of 1737 ended the theatrical career of Henry Fielding, whose comedies had come under constant fire from the authorities for their satire on the government. Fielding's comic talents were perforce directed to the novel, the form in which he parodied the sentiment and the morality of Richardson's *Pamela*—in his *Shamela* and *Joseph Andrews* (1742)—as brilliantly as he had earlier burlesqued the rant of heroic tragedy in *Tom Thumb* (1730).

Comedy of the sort that ridicules the follies and vices of society to the end of laughing them out of countenance entered the English novel with Fielding. His statement in *Joseph Andrews* concerning the function of satire is squarely in the Neoclassic tradition of comedy as a corrective of manners and mores: the satirist holds

the glass to thousands in their closets, that they may contemplate their deformity, and endeavour to reduce it, and thus by suffering private mortification may avoid public shame.

Fielding's scenes of contemporary life display the same power of social criticism as that which distinguishes the engravings of his great fellow artist William Hogarth,

whose "Marriage à la Mode" (1745) depicts the vacuity and the casual wantonness of the fashionable world that Fielding treats of in the final books of *Tom Jones*. Hogarth's other series, such as "A Rake's Progress" (1735) or "A Harlot's Progress" (1732), also make a didactic point about the wages of sin, using realistic details heightened with grotesquerie to expose human frailty and its sinister consequences. The grotesque is a recurrent feature of the satiric tradition in England, where comedy serves social criticism. Artists such as Hogarth and Thomas Rowlandson worked in the tradition of Jonson and the Restoration dramatists in the preceding century.

The novel, with its larger scope for varied characters, scenes, and incidents, rather than the drama, afforded the 19th-century artist in comedy a literary form adequate to his role as social critic. The spectacle of man and his society is regularly presented by the 19th-century novelist in comedic terms, as in *Vanity Fair* (1848), by William Makepeace Thackeray or the *Comédie humaine* (1842–55) of Honoré de Balzac, and with the novels of Jane Austen, Anthony Trollope, Charles Dickens, and George Meredith.

20th-century tragicomedy. The best that the comic stage had to offer in the late 19th century lay in the domain of farce. The masters of this form were French, but it flourished in England as well; what the farces of Eugène Labiche and Georges Feydeau and the operettas of Jacques Offenbach were to the Parisian stage the farces of W.S. Gilbert and the young Arthur Wing Pinero and the operettas that Gilbert wrote in collaboration with Arthur Sullivan were to the London stage. As concerns comedy, the situation in England improved at the end of the century, when Oscar Wilde and George Bernard Shaw turned their talents to it. Wilde's *Importance of Being Earnest* (1895) is farce raised to the level of high comic burlesque. Shaw's choice of the comic form was inevitable, given his determination that the contemporary English stage should deal seriously and responsibly with the issues that were of crucial importance to contemporary English life. Serious subjects could not be resolved by means of the dramatic clichés of Victorian melodrama. Rather, the prevailing stereotypes concerning the nature of honour, courage, wisdom, and virtue were to be subjected to a hail of paradox, to the end of making evident their inner emptiness or the contradictions they concealed.

Shaw dealt with what, in the preface to *Major Barbara* (1905), he called "the tragi-comic irony of the conflict between real life and the romantic imagination," and his use of the word tragicomic is a sign of the times. The striking feature of modern art, according to the German novelist Thomas Mann, was that it had ceased to recognize the categories of tragic and comic or the dramatic classifications of tragedy and comedy but saw life as tragicomedy. The sense that tragicomedy is the only adequate dramatic form for projecting the unreconciled ironies of modern life mounted through the closing decades of the 19th century. Ibsen had termed *The Wild Duck* (published 1884) a tragicomedy; it was an appropriate designation for this bitter play about a young man blissfully ignorant of the lies on which he and his family have built their happy life until an outsider who is committed to an ideal of absolute truth exposes all their guilty secrets with disastrous results. The plays of the Russian writer Anton Chekhov, with their touching and often quite humorous figures leading lives of quiet desperation, reflect precisely that mixture of inarticulate joy and dull pain that is the essence of the tragicomic view of life.

A dramatist such as August Strindberg produces a kind of tragicomedy peculiarly his own, one that takes the form of bourgeois tragedy; it lacerates its principals until they become a parody of themselves. Strindberg's *Dance of Death* (1901), with its cruelty and pain dispensed with robust pleasure by a fiercely battling husband and wife, is a significant model of the grotesque in the modern theatre; it is reflected in such mid-20th-century examples of what came to be called black comedy as Eugène Ionesco's *Victims of Duty* (1953) and Edward Albee's *Who's Afraid of Virginia Woolf?* (1962). Almost equally influential as a turn-of-the-century master of the grotesque is Frank Wedekind, whose

French farce

Censorship

Earth Spirit (1895) and its sequel, *Pandora's Box* (written 1892–1901), though both are termed tragedies by their author, are as much burlesques of tragedy as *The Dance of Death*. Their grotesquerie consists chiefly in their disturbing combination of innocence and depravity, of farce and horror, of passionate fervour issuing in ludicrous incident that turns deadly. Wedekind's celebration of primitive sexuality and the varied ways in which it manifests itself in an oversophisticated civilization distorts the tragic form to achieve its own grotesque beauty and power.

The great artist of the grotesque and of tragicomedy in the 20th century is the Italian Luigi Pirandello. His drama is explicitly addressed to the contradictoriness of experience: appearances collide and cancel out each other; the quest of the absolute issues in a mind-reeling relativism; infinite spiritual yearnings are brought up hard against finite physical limits; rational purpose is undermined by irrational impulse; and with the longing for permanence in the midst of change comes the ironic awareness that changelessness means death. Stated thus, Pirandello's themes sound almost forbiddingly intellectual, but one of his aims was to convert intellect into passion. Pirandello's characters suffer from intellectual dilemmas that give rise to mental and emotional distress of the most anguished kind, but their sufferings are placed in a satiric frame. The incongruities that the characters are furiously seeking to reconcile attest to the comic aspect of this drama, but there is nothing in it of the traditional movement of comedy, from a state of illusion into the full light of reality. Pirandello's characters dwell amid ambiguities and equivocations that those who are wise in the tragicomic nature of life will accept without close inquiry. The logic of comedy implies that illusions exist to be dispelled; once they are dispelled, everyone will be better off. The logic of Pirandello's tragicomedy demonstrates that illusions make life bearable; to destroy them is to destroy the basis for any possible happiness.

The role of
illusions in
Pirandello's
tragi-
comedy

The absurd. In their highly individual ways, both Samuel Beckett and Ionesco have employed the forms of comedy—from tragicomedy to farce—to convey the vision of an exhausted civilization and a chaotic world. The very endurance of life amid the grotesque circumstances that obtain in Beckett's plays is at once a tribute to the human power of carrying on to the end and an ironic reflection on the absurdity of doing so. Beckett's plays close in an uneasy silence that is the more disquieting because of the uncertainty as to just what it conceals: whether it masks sinister forces ready to spring or is the expression of a universal indifference or issues out of nothing at all.

Silence seldom reigns in the theatre of Ionesco, which rings with voices raised in a usually mindless clamour. Some of Ionesco's most telling comic effects have come from his use of dialogue overflowing with clichés and non sequiturs, which make it clear that the characters do not have their minds on what they are saying and, indeed, do not have their minds on anything at all. What they say is often at grotesque variance with what they do. Beneath the moral platitudes lurks violence, which is never far from the surface in Ionesco's plays, and the violence tells what happens to societies in which words and deeds have become fatally disjunct. Ionesco's comic sense is evident as well in his depiction of human beings as automata, their movements decreed by forces they have never questioned or sought to understand. There is something undeniably farcical in Ionesco's spectacles of human regimentation, of men and women at the mercy of things (e.g., the stage full of chairs in *The Chairs* or the growing corpse in *Amédée*); the comic quality here is one that Bergson would have appreciated. But the comic in Ionesco's most serious work, as in so much of the contemporary theatre, has ominous implications that give to it a distinctly grotesque aspect. In Ionesco's *Victims of Duty* and *The Killer* (1959), as in the works of his Swiss counterparts—*Der Besuch der alten Dame* (performed 1956; *The Visit*, 1958) and *The Physicists* (1962), by Friedrich Dürrenmatt, and *The Firebugs* (1958), by Max Frisch—the grotesquerie of the tragicomic vision delineates a world in which the humane virtues are dying, and casual violence is the order of the day.

Ionesco's
comic
sense

The radical reassessment of the human image that the

20th century has witnessed is reflected in the novel as well as in drama. Previous assumptions about the rational and divine aspects of man have been increasingly called into question by the evidences of man's irrationality, his sheer animality. These are qualities of human nature that writers of previous ages (Swift, for example) have always recognized, but hitherto they have been typically viewed as dark possibilities that could overtake humanity if the rule of reason did not prevail. It is only in the mid-20th century that the savage and the irrational have come to be viewed as part of the normative condition of humanity rather than as tragic aberrations from it. The savage and the irrational amount to grotesque parodies of human possibility, ideally conceived. Thus it is that 20th-century novelists as well as dramatists have recognized the tragicomic nature of the contemporary human image and predicament, and the principal mode of representing both is the grotesque. This may take various forms: the apocalyptic nightmare of tyranny and terror in Kafka's novels *The Trial* (1925) and *The Castle* (1926); the tragic farce in terms of which the Austrian novelist Robert Musil describes the slow collapse of a society into anarchy and chaos, in *The Man Without Qualities* (1930–43); the brilliant irony whereby Thomas Mann represents the hero as a confidence man in *The Confessions of Felix Krull* (1954); the grimly parodic account of Germany's descent into madness in Günter Grass's novel *The Tin Drum* (1959). The English novel contains a rich vein of the comic grotesque that extends at least back to Dickens and Thackeray and persisted in the 20th century in such varied novels as Evelyn Waugh's *Decline and Fall* (1928), Angus Wilson's *Anglo-Saxon Attitudes* (1956), and Kingsley Amis' *Lucky Jim* (1954). What novelists such as these have in common is the often disturbing combination of hilarity and desperation. It has its parallel in a number of American novels—John Barth's *Giles Goat-Boy* (1966), Kurt Vonnegut, Jr.'s *Slaughterhouse Five* (1969)—in which shrill farce is the medium for grim satire. And the grotesque is a prominent feature of modern poetry, as in some of the "Songs and Other Musical Pieces" of W.H. Auden.

THE COMIC IN OTHER ARTS

The visual arts. The increasing use of the affairs of common life as the subject matter of dramatic comedy through the Middle Ages and the Renaissance is also seen in painting of that time. Scenes from medieval mystery cycles, such as the comic episodes involving Noah's stubborn wife, have counterparts in medieval pictures in the glimpses of everyday realities that are caught through the windows or down the road from the sites where the great spiritual mysteries are in progress: the angel Gabriel may appear to the Virgin in the foreground, while a man is chopping wood in the yard outside. Medieval artists had never neglected the labours and the pleasures of the mundane world, but the treatment of them is often literally marginal, as in the depiction of men and women at work or play in the ornamental borders of an illuminated manuscript page. The seasonal round of life, with its cycle of plowing, sowing, mowing, and reaping interspersed with hawking, hunting, feasts, and weddings (the cycle of life, indeed, which comedy itself celebrates), is depicted in series after series of exquisite miniatures, such as those in the *Très Riches Heures du Duc de Berry*. By the mid-16th century, however, in Pieter Bruegel's famous painting "Landscape with the Fall of Icarus," mundane reality has taken over the foreground; the plowman tills the soil, and the shepherd attends his flock, while, unnoticed by both, the legs of Icarus disappear inconspicuously into the sea. Bruegel is not a comic artist, but his art bears witness to what all great comic art celebrates: the basic rhythm of life. "Peasant Wedding" and "Peasant Dance" endow their heavy men and women with an awkward grace and dignity that bear comparison with Shakespeare's treatment of his comic characters. Paintings like Bruegel's "Children's Games" and his "Fight Between Carnival and Lent" are joyous representations of human energy. The series of "The Labours of the Months"—"Hunters in the Snow" for January, "Haymaking" for July, "Harvesters" for August, "Return of the Herd" for November—give

Bruegel's
mastery
of the
grotesque

pictorial treatment to a favourite subject of the medieval miniaturists. Finally, allusion must be made to Bruegel's mastery of the grotesque, notably in "The Triumph of Death" and in the "Dulle Griet," in which demons swarm over a devastated landscape.

It is through the art of caricature that the spirit of comedy enters most directly into painting. The style derives from the portraits with ludicrously exaggerated features made by the Carracci, an Italian family of artists, early in the 17th century (Italian *caricare*, "to overload"). In defiance of the theory of ideal beauty, these portraits emphasized the features that made one man different from another. This method of character portrayal—the singling out of one distinctive feature and emphasizing it over all others—is not unlike the practice of characterizing the personages of the comic stage by means of some predominant humour, which Ben Jonson was developing at about the same time in the London theatre. The use of exaggeration for comic effect was as evident to painters as it was to dramatists. Its usefulness as a means of social and political satire is fully recognized by Hogarth. Hogarth's counterpart in mid-19th-century Paris was Honoré Daumier. His caricatures portray a human comedy as richly detailed and as shrewdly observed as the one portrayed in fiction by his contemporary Balzac. But Daumier's sense of the comic goes beyond caricature; his numerous treatments of scenes from Molière's plays and, most notably, his drawings and canvases of Don Quixote and Sancho Panza attest to the pathos that can lie beneath the comic mask.

Modern art has abstracted elements of comedy to aid it in the representation of a reality in which the mechanical is threatening to win out over the human. Bergson's contention that the essence of comedy consists of something mechanical encrusted on the living may be said to have achieved a grotesque apotheosis in the French Dadaist Marcel Duchamp's painting "Bride" (1912), in which the female figure has been reduced to an elaborate piece of plumbing. The highly individual Swiss Expressionist Paul Klee's pen-and-ink drawing tinted with watercolour and titled "Twittering Machine" (1922) represents an ingenious device for imitating the sound of birds. The delicacy of the drawing contrasts with the sinister implications of the mechanism, which, innocent though it may appear at first glance, is almost certainly a trap.

The grotesque is a constant stylistic feature of the artist's representation of reality in its brutalized or mechanized aspects. The carnival masks worn by the figures in the painting "Intrigue" (1890), by the Belgian James Ensor, make manifest the depravity and the obscenity that lurk beneath the surface of conventional appearances; Ensor's paintings make much the same point about the persistence of the primitive and the savage into modern life as Wedekind's plays were to do a few years later. German artists after World War I invoked the grotesque with particular power, depicting the inhuman forces that bear upon the individual, as in George Grosz's savage cartoon titled "Germany, a Winter's Tale" (1918), in which the puppet-like average citizen sits at table surrounded by militarist, capitalist, fatuous clergyman and all the violent and dissolute forces of a decadent society. The mutilated humanity in Max Beckmann's "Dream" (1921) and "Departure" (1932–33) is a further testament to human viciousness, 20th-century variety.

Rather more explicitly comic is the element of fantasy in modern paintings, in which seemingly unrelated objects are brought together in a fine incongruity, as in the French primitive Henri Rousseau's famous "Dream" (1910), with its nude woman reclining on a red-velvet sofa amid the flora and fauna of a lush and exotic jungle. The disparate figures that float (in defiance of all the laws of gravity) through the paintings of the Russian Surrealist Marc Chagall are individually set forth in a nimbus of memory and in the landscape of dream. But fantasy can take on a grotesquerie of its own, as in some of Chagall's work, such as the painting "I and The Village" (1911).

The purest expression of the comic in modern painting must surely be Henri Matisse's "Joy of Life" (1905–06), a picture that might be taken as a visual expression of the precept that the rhythm of comedy is the basic rhythm of

life. But Matisse's painting was not to be the last word on the subject: "Joy of Life" produced, as a counterstatement, Pablo Picasso's "Demoiselles d'Avignon" (1906–07), in which the daughters of joy, in their grim and aggressive physical tension, stand as a cruel parody of the delight in the senses that Matisse's picture celebrates. "Les Demoiselles d'Avignon" and such a later Picasso masterpiece as the "Three Dancers" (1925) suggest that, for the visual as well as the literary artist of the 20th century, the joy of life tends to issue in grotesque shapes.

Music. Given the wide range of imitative sounds of which musical instruments and the human voice are capable, comic effects are readily available to the composer who wants to use them. At the simplest level, these may amount to nothing more than humorous adjuncts to a larger composition, such as the loud noise with which the 18th-century Austrian composer Joseph Haydn surprises his listeners in *Symphony No. 94* or the sound of the ticking clock in *No. 101*. The scherzo, which Ludwig van Beethoven introduced into symphonic music in the early 19th century, may be said to have incorporated in it a musical joke but one of a highly abstract kind; its nervous jocularity provides a contrast and a commentary (both heavy with irony) on the surrounding splendour. A century after Beethoven, the jocularity grew more desperate and the irony more profound in the grim humour that rises out of the grotesque scherzos of Gustav Mahler. A more sustained and a more explicit musical exposition of comic themes and attitudes comes when a composer draws his inspiration directly from a work of comic literature, as Richard Strauss does in his orchestral variations based on *Don Quixote* and on the merry pranks of *Till Eulenspiegel*.

It is, however, opera that provides the fullest form for comedy to express itself in music, and some of the most notable achievements of comic art have been conceived for the operatic stage. High on any list of comic masterpieces must come the four principal operas of Mozart: *The Marriage of Figaro* (1786), *Don Giovanni* (1787), *Così fan tutte* (1790), and *The Magic Flute* (1791), and there are countless others worthy of mention. Operatic comedy has an advantage over comedy in the spoken theatre in its ability to impose a coherent form on the complexities of feeling and action that are often of the essence in comedy. The complex feeling experienced by different characters must be presented in spoken comedy seriatim; operatic comedy can present them simultaneously. When three or four characters talk simultaneously in the spoken theatre, the result is an incoherent babble. But the voices of three or four or even more characters can be blended together in an operatic ensemble, and, while most of the words may be lost, the vocal lines will serve to identify the individual characters and the general nature of the emotions they are expressing. The complexities of action in the spoken theatre are the chief source of the comic effect, which increases as the confusion mounts; such complexities of action operate to the same comic end in opera but here with the added ingredient of music, which provides an overarching design of great formal coherence. In the music, all is manifestly ordered and harmonious, while the events of the plot appear random and chaotic; the contrast between the movement of the plot and the musical progression provides a Mozart or a Rossini with some of his wittiest and most graceful comic effects. Finally, it should be noted that operatic comedy can probe psychological and emotional depths of character that spoken comedy would scarcely attempt. The Countess in Mozart's *Figaro* is a very much more moving figure than she is in Beaumarchais' play; the Elvira of *Don Giovanni* exhibits a fine extravagance that is little more than suggested in Molière's comedy.

Television and cinema. When comedy is dependent on the favour of a large part of the public, as reflected in box-office receipts or the purchase of a television sponsor's product, it seldom achieves a high level of art. There is nothing innocent about laughter at the whims and inconsistencies of humankind, and radio and television and film producers have always been wary of offending their audiences with it. On radio and television, the laughter is

The form
imposed by
music

Fantasy in
modern
paintings

usually self-directed (as in the performances of comedians such as Jack Benny or Red Skelton), or it is safely contained within the genial confines of a family situation (e.g., the "Fibber McGee and Molly" radio show or "I Love Lucy" on television). Much the same attitude has obtained with regard to comedy in the theatre in the United States. Satire has seldom succeeded on Broadway, which instead has offered pleasant plays about the humorous behaviour of basically nice people, such as the eccentric family in George S. Kaufmann and Moss Hart's *You Can't Take It with You* (1936) or the lovable head of the household in Howard Lindsay and Russel Crouse's *Life with Father* (1939) or the indefatigable Dolly Levi in Thornton Wilder's *Matchmaker* (1954) and in her later reincarnation in the musical *Hello, Dolly!*

The American public and comedy

The American public has never been quite comfortable in the presence of comedy. The calculated ridicule and the relentless exposure often seem cruel or unfair to a democratic public. If all men are created equal, then it ill becomes anyone to laugh at the follies of his fellows, especially when they are follies that are likely to be shared, given the common background of social opportunity and experience of the general public. There is an insecurity in the mass audience that is not compatible with the high self-assurance of comedy as it judges between the wise and the foolish of the world. The critical spirit of comedy has never been welcome in American literature; in both fiction and drama, humour, not comedy, has raised the laughter. American literature can boast an honorable tradition of humorists, from Mark Twain to James Thurber, but has produced no genuinely comic writer. As American social and moral tenets were subjected to increasing critical scrutiny from the late 1950s onward, however, there were some striking achievements in comedy in various media: Edward Albee's *American Dream* (1961) and *Who's Afraid of Virginia Woolf?* (1962), on the stage; novels such as those of Saul Bellow and Joseph Heller's *Catch-22* (1961); and films such as *Dr. Strangelove* (1964).

This last example is remarkable, because comedy in the medium of film, in America, had been conceived as entertainment and not much more. This is not to say that American film comedies lacked style. The best of them always displayed verve and poise and a thoroughly professional knowledge of how to amuse the public without troubling it. Their shortcoming has always been that the amusement they provide lacks resonance.

If films have seldom explored comedy with great profundity, they have, nonetheless, produced it in great variety. There have been comedies of high sophistication, the work of directors such as Ernst Lubitsch, George Cukor, Frank Capra, Joseph L. Mankiewicz, and Billy Wilder and of actors and actresses such as Greta Garbo (in Lubitsch's *Ninotchka*, 1939), Katharine Hepburn and Cary Grant (in Cukor's *Philadelphia Story*, 1940), Bette Davis (in Mankiewicz's *All About Eve*, 1950), Clark Gable and Claudette Colbert (in Capra's *It Happened One Night*, 1934), Gary Cooper and Jean Arthur (in Capra's *Mr. Deeds Goes to Town*, 1936), and Marilyn Monroe and Jack Lemmon (in Billy Wilder's *Some Like It Hot*, 1959). There have been comedies with music, built around the talents of singers and dancers such as Ruby Keeler and Dick Powell and Ginger Rogers and Fred Astaire; there are the classic farces of Charlie Chaplin and Buster Keaton and, later, of W.C. Fields and the Marx Brothers and Laurel and Hardy; and there is a vast, undistinguished field of comedies dealing with the humours of domestic life. The varieties of comedy in Hollywood films have always been replicas of those on the New York stage; as often as not, they were products of the same talents: in the 1930s, of dramatists such as Philip Barry or S.N. Behrman and composers such as Cole Porter, Richard Rodgers, and Irving Berlin; in the 1960s, of the dramatist Neil Simon and the composer Burt Bacharach.

Comedies produced by European film makers

European film makers, with an older and more intellectual tradition of comedy available to them, produced comedies of more considerable stature. Among French directors, Jean Renoir, in his *The Rules of the Game* (1939), conveyed a moving human drama and a profoundly serious vision of French life on the eve of World

War II in a form, deriving from the theatre, that blends the comic and the tragic. His disciple François Truffaut, in *Jules and Jim* (1961), directed a witty and tender but utterly clear-sighted account of how gaiety and love turn deadly. Though not generally regarded as a comic artist, the Swedish film maker Ingmar Bergman produced a masterpiece of film comedy in *Smiles of a Summer Night* (1955), a wise, wry account of the indignities that must sometimes be endured by those who have exaggerated notions of their wisdom or virtue. The films of the Italian director and writer Federico Fellini represent a comic vision worthy of Pirandello. *La strada* (1954), with its Chaplinesque waif (played by Fellini's wife, Giulietta Masina) as central figure, is a disturbing compound of pathos and brutality. Comedy's affirmation of the will to go on living has had no finer portrayal than in Giulietta Masina's performance in the closing scene of *Nights of Cabiria* (1956). *La dolce vita* (1960) is a luridly satiric vision of modern decadence, where ideals are travestied by reality, and everything is illusion and disillusionment; the vision is carried to even more bizarre lengths in Fellini's *Satyricon* (1969), in which the decadence of the modern world is grotesquely mirrored in the ancient one. *8½* (1963) and *Juliet of the Spirits* (1965) are Fellini's most brilliantly inventive films, but their technical exuberance is controlled by a profoundly serious comic purpose. The principals in both films are seeking—through the phantasmagoria of their past and present, of their dreams and their delusions, all of which seem hopelessly mixed with their real aspirations—to know themselves. (C.H.Ho.)

Tragedy

Although the word tragedy is often used loosely to describe any sort of disaster or misfortune, it more precisely refers to a work of art, usually a play or novel, that probes with high seriousness questions concerning the role of man in the universe. The Greeks of Attica, the ancient state whose chief city was Athens, first used the word in the 5th century BC to describe a specific kind of play, which was presented at festivals in Greece. Sponsored by the local governments, these plays were attended by the entire community, a small admission fee being provided by the state for those who could not afford it themselves. The atmosphere surrounding the performances was more like that of a religious ceremony than entertainment. There were altars to the gods, with priests in attendance, and the subjects of the tragedies were the misfortunes of the heroes of legend, religious myth, and history. Most of the material was derived from the works of Homer and was common knowledge in the Greek communities. So powerful were the achievements of the three greatest Greek dramatists—Aeschylus (525–456 BC), Sophocles (c. 496–406 BC), and Euripides (c. 480–406 BC)—that the word they first used for their plays survived and came to describe a literary genre that, in spite of many transformations and lapses, has proved its viability through 25 centuries.

Historically, tragedy of a high order has been created in only four periods and locales: Attica, in Greece, in the 5th century BC; England in the reigns of Elizabeth I and James I, from 1558 to 1625; 17th-century France; and Europe and America during the second half of the 19th century and the first half of the 20th. Each period saw the development of a special orientation and emphasis, a characteristic style of theatre. In the modern period, roughly from the middle of the 19th century, the idea of tragedy found embodiment in the collateral form of the novel.

This section focusses primarily on the development of tragedy as a literary genre. Further information on the relationship of tragedy to other types of drama will be found in the section above on *Dramatic literature*. The role of tragedy in the growth of theatre is discussed in THEATRE, THE HISTORY OF WESTERN.

DEVELOPMENT

Origins in Greece. The questions of how and why tragedy came into being and of the bearing of its origins on its development in subsequent ages and cultures have been

Four great periods of tragedy

investigated by historians, philologists, archaeologists, and anthropologists with results that are suggestive but conjectural. Even the etymology of the word tragedy is far from established. The most generally accepted source is the Greek *tragōidia*, or "goat-song," from *tragos* ("goat") and *aeidein* ("to sing"). The word could have referred either to the prize, a goat, that was awarded to the dramatists whose plays won the earliest competitions or to the dress (goat skins) of the performers, or to the goat that was sacrificed in the primitive rituals from which tragedy developed.

In these communal celebrations, a choric dance may have been the first formal element and perhaps for centuries was the principal element. A speaker was later introduced into the ritual, in all likelihood as an extension of the role of the priest, and dialogue was established between him and the dancers, who became the chorus in the Athenian drama. Aeschylus is usually regarded as the one who, realizing the dramatic possibilities of the dialogic, first added a second speaker and thus invented the form of tragedy. That so sophisticated a form could have been fully developed by a single artist, however, is scarcely credible. Hundreds of early tragedies have been lost, including some by Aeschylus himself. Of some 90 plays attributed to him, only seven have survived.

Four Dionysia, or feasts of the Greek God Dionysus, were held annually in Athens. Since Dionysus once held place as the god of vegetation and the vine, and the goat was believed sacred to him, it has been conjectured that tragedy originated in fertility feasts to commemorate the harvest and the vintage and the associated ideas of the death and renewal of life. The purpose of such rituals is to exercise some influence over these vital forces. Whatever the original religious connections of tragedy may have been, two elements have never entirely been lost: (1) its high seriousness, befitting matters in which survival is at issue and (2) its involvement of the entire community in matters of ultimate and common concern. When either of these elements diminishes, when the form is overmixed with satiric, comic, or sentimental elements, or when the theatre of concern succumbs to the theatre of entertainment, then tragedy falls from its high estate and is on its way to becoming something else.

As the Greeks developed it, the tragic form, more than any other, raised questions about man's existence. Why must man suffer? Why must man be forever torn between the seeming irreconcilables of good and evil, freedom and necessity, truth and deceit? Are the causes of his suffering outside himself, in blind chance, in the evil designs of others, in the malice of the gods? Are its causes within him, and does he bring suffering upon himself through arrogance, infatuation, or the tendency to overreach himself? Why is justice so elusive?

Aeschylus: the first great tragedian. It is this last question that Aeschylus asks most insistently in his two most famous works, the *Oresteia* (a trilogy comprising *Agamemnon*, *Choephoroi*, and *Eumenides*) and *Prometheus Bound* (the first part of a trilogy of which the last two parts have been lost): is it right that Orestes, a young man in no way responsible for his situation, should be commanded by a god, in the name of justice, to avenge his father by murdering his mother? Is there no other way out of his dilemma than through the ancient code of blood revenge, which will only compound the dilemma? Again: was it right that in befriending mankind with the gifts of fire and the arts, Prometheus should offend the presiding god Zeus and himself be horribly punished? Aeschylus opened questions whose answers in the Homeric stories had been taken for granted. In Homer, Orestes' patricide is regarded as an act of filial piety, and Prometheus' punishment is merely the inevitable consequence of defying the reigning deity. All of the materials of tragedy, all of its cruelty, loss, and suffering, are present in Homer and the ancient myths but are dealt with as absolutes—self-sufficient and without the questioning spirit that was necessary to raise them to the level of tragedy. It remained for Aeschylus and his fellow tragedians first to treat these "absolutes" critically and creatively in sustained dramatic form. They were true explorers of the human spirit.

In addition to their remarkable probing into the nature

of existence, their achievements included a degree of psychological insight for which they are not generally given credit. Though such praise is usually reserved for Shakespeare and the moderns, the Athenian dramatists conveyed a vivid sense of the living reality of their characters' experience: of what it felt like to be caught, like Orestes, in desperately conflicting loyalties or to be subjected, like Prometheus, to prolonged and unjust punishment. The mood of the audience as it witnessed the acting out of these climactic experiences has been described as one of impassioned contemplation. From their myths and epics and from their history in the 6th century, the people of Athens learned that they could extend an empire and lay the foundations of a great culture. From their tragedies of the 5th century, they learned who they were, something of the possibilities and limitations of the spirit, and of what it meant, not merely what it felt like, to be alive in a world both beautiful and terrible.

Aeschylus has been called the most theological of the Greek tragedians. His *Prometheus* has been compared to the Book of Job of the Bible both in its structure (*i.e.*, the immobilized heroic figure maintaining his cause in dialogues with visitors) and in its preoccupation with the problem of suffering at the hands of a seemingly unjust deity. Aeschylus tended to resolve the dramatic problem into some degree of harmony, as scattered evidence suggests he did in the last two parts of the *Promethiad* and as he certainly did in the conclusion of the *Oresteia*. This tendency would conceivably lead him out of the realm of tragedy and into religious assurance. But his harmonies are never complete. In his plays evil is inescapable, loss is irretrievable, suffering is inevitable. What the plays say positively is that man can learn through suffering. The chorus in *Agamemnon*, the first play of the *Oresteia*, says this twice. The capacity to learn through suffering is a distinguishing characteristic of the tragic hero, preeminently of the Greek tragic hero. He has not merely courage, tenacity, and endurance but also the ability to grow, by means of these qualities, into an understanding of himself, of his fellows, and of the conditions of existence. Suffering, says Aeschylus, need not be embittering but can be a source of knowledge. The moral force of his plays and those of his fellow tragedians can hardly be exaggerated. They were shaping agents in the Greek notion of education. It has been said that from Homer the Greeks learned how to be good Greeks; from the tragedies they learned an enlarged humanity. If it cannot be proved that Aeschylus "invented" tragedy, it is clear that he at least set its tone and established a model that is still operative. Even in the 20th century, the *Oresteia* has been acclaimed as the greatest spiritual work of man, and dramatists such as T.S. Eliot, in *The Family Reunion* (1939), and Jean-Paul Sartre, in *The Flies* (1943), found modern relevance in its archetypal characters, situations, and themes.

Sophocles: the purest artist. Sophocles' life spanned almost the whole of the 5th century. He is said to have written his last play, *Oedipus at Colonus*, at the age of 90. Only seven of his plays, of some 125 attributed to him, survive. He won the prize in the tragic competitions 20 times and never placed lower than second.

Sophocles has been called the great mediating figure between Aeschylus and Euripides. Of the three, it might be said that Aeschylus tended to resolve tragic tensions into higher truth, to look beyond, or above, tragedy; that Euripides' irony and bitterness led him the other way to fix on the disintegration of the individual; and that Sophocles, who is often called the "purest" artist of the three, was truest to the actual state of human experience. Unlike the others, Sophocles seems never to insinuate himself into his characters or situations, never to manipulate them into preconceived patterns. He sets them free on a course seemingly of their own choosing. He neither preaches nor rails. If life is hard and often destructive, the question Sophocles asks is not how did this come to be or why did such a misfortune have to happen but rather, given the circumstances, how must a man conduct himself, how should he act, what must he do?

His greatest play, *Oedipus the King*, may serve as a model of his total dramatic achievement. Embodied in it,

Aeschylus' view of the problem of suffering

Two vital elements

and suggested with extraordinary dramatic tact, are all the basic questions of tragedy, which are presented in such a way as almost to define the form itself. It is not surprising that Aristotle, a century later, analyzed it for his definition of tragedy in the *Poetics*. It is the nuclear Greek tragedy, setting the norm in a way that cannot be claimed for any other work, not even the *Oresteia*.

Sophocles' emphasis on action

In *Oedipus*, as in Sophocles' other plays, the chorus is much less prominent than in Aeschylus' works. The action is swifter and more highly articulated; the dialogue is sharper, more staccato, and bears more of the meaning of the play. Though much has been made of the influence of fate on the action of the play, later critics emphasize the freedom with which Oedipus acts throughout. Even before the action of the play begins, the oracle's prediction that Oedipus was doomed to kill his father and marry his mother had long since come true, though he did not realize it. Though he was fated, he was also free throughout the course of the play—free to make decision after decision, to carry out his freely purposed action to its completion. In him, Sophocles achieved one of the enduring definitions of the tragic hero—that of a man for whom the liberation of the self is a necessity. The action of the play, the purpose of which is to discover the murderer of Oedipus' father and thereby to free the city from its curse, leads inevitably to Oedipus' suffering—the loss of his wife, his kingdom, his sight. The messenger who reports Oedipus' self-blinding might well have summarized the play with "All ills that there are names for, all are here." And the chorus' final summation deepens the note of despair: "Count no man happy," they say in essence, "until he is dead."

But these were not Sophocles' ultimate verdicts. The action is so presented that the final impression is not of human helplessness at the hands of maligning gods nor of man as the pawn of fate. Steering his own course, with great courage, Oedipus has ferreted out the truth of his identity and administered his own punishment, and, in his suffering, learned a new humanity. The final impression of the *Oedipus*, far from being one of unmixed evil and nihilism, is of massive integrity, powerful will, and unanimous acceptance of a horribly altered existence.

Some 50 years later, Sophocles wrote a sequel to *Oedipus the King*. In *Oedipus at Colonus*, the old Oedipus, further schooled in suffering, is seen during his last day on earth. He is still the same Oedipus in many ways: hot-tempered, hating his enemies, contentious. Though he admits his "pollution" in the murder of his father and the marriage to his mother, he denies that he had sinned, since he had done both deeds unwittingly. Throughout the play, the theme of which has been described as the "heroization" of Oedipus, he grows steadily in nobility and awesomeness. Finally, sensing the approach of the end, he leaves the scene, to be elevated in death to a demigod, as the messenger describes the miraculous event. In such manner Sophocles leads his tragedy toward an ultimate assertion of values. His position has been described as "heroic humanism," as making a statement of belief in the human capacity to transcend evils, within and without, by means of the human condition itself.

Tragedy must maintain a balance between the higher optimisms of religion or philosophy, or any other beliefs that tend to explain away the enigmas and afflictions of existence, on the one hand, and the pessimism that would reject the whole human experiment as valueless and futile on the other. Thus the opposite of tragedy is not comedy but the literature of cynicism and despair, and the opposite of the tragic artist's stance, which is one of compassion and involvement, is that of the detached and cynical ironist.

Euripides: the dark tragedian. The tragedies of Euripides test the Sophoclean norm in this direction. His plays present in gruelling detail the wreck of human lives under the stresses that the gods often seem willfully to place upon them. Or, if the gods are not willfully involved through jealousy or spite, they sit idly by while man wrecks himself through passion or heedlessness. No Euripidean hero approaches Oedipus in stature. The margin of freedom is narrower, and the question of justice, so central and absolute an ideal for Aeschylus, becomes a subject for

Euripides' view of justice

irony. In *Hippolytus*, for example, the goddess Artemis never thinks of justice as she takes revenge on the young Hippolytus for neglecting her worship; she acts solely out of personal spite. In *Medea*, Medea's revenge on Jason through the slaughter of their children is so hideously unjust as to mock the very question. In the *Bacchae*, when the frenzied Agave tears her son, Pentheus, to pieces and marches into town with his head on a pike, the god Dionysus, who had engineered the situation, says merely that Pentheus should not have scorned him. The Euripidean gods, in short, cannot be appealed to in the name of justice. Euripides' tendency toward moral neutrality, his cool tacking between sides (e.g., between Pentheus versus Dionysus and the bacchantes) leave the audience virtually unable to make a moral decision. In Aeschylus' *Eumenides* (the last play of the *Oresteia*), the morals of the gods improve. Athena is there, on the stage, helping to solve the problem of justice. In Sophocles, while the gods are distant, their moral governance is not questioned. *Oedipus* ends as if with a mighty "So be it." In Euripides, the gods are destructive, wreaking their capricious wills on defenseless man. Aristotle called Euripides the most tragic of the three dramatists; surely his depiction of the arena of human life is the grimmest.

Many qualities, however, keep his tragedies from becoming literature of protest, of cynicism, or of despair. He reveals profound psychological insight, as in the delineation of such antipodal characters as Jason and Medea, or of the forces, often subconscious, at work in the group frenzy of the *Bacchae*. His Bacchic odes reveal remarkable lyric power. And he has a deep sense of human values, however external and self-conscious. Medea, even in the fury of her hatred for Jason and her lust for revenge, must steel herself to the murder of her children, realizing the evil of what she is about to do. In this realization, Euripides suggests a saving hope: here is a great nature gone wrong—but still a great nature.

Later Greek drama. After Euripides, Greek drama reveals little that is significant to the history of tragedy. Performances were given during the remainder of the pre-Christian era in theatres throughout the Mediterranean world, but, with the decline of Athens as a city-state, the tradition of tragedy eroded. As external affairs deteriorated, the high idealism, the exalted sense of human capacities depicted in tragedy at its height yielded more and more to the complaints of the skeptics. The Euripidean assault on the gods ended in the debasement of the original lofty conceptions. A 20th-century British classical scholar, Gilbert Murray, used the phrase "the failure of nerve" to describe the late Greek world. It may, indeed, provide a clue to what happened. On the other hand, according to the 19th-century German philosopher Friedrich Nietzsche, in *The Birth of Tragedy* (1872), a quite different influence may have spelled the end of Greek tragedy: the so-called Socratic optimism, the notion underlying the dialogues of Plato that man could "know himself" through the exercise of his reason in patient, careful dialectic—a notion that diverted questions of man's existence away from drama and into philosophy. In any case, the balance for tragedy was upset, and the theatre of Aeschylus, Sophocles, and Euripides gave way to what seems to have been a theatre of diatribe, spectacle, and entertainment.

The long hiatus. The Roman world failed to revive tragedy. Seneca (4 BC–AD 65) wrote at least eight tragedies, mostly adaptations of Greek materials, such as the stories of Oedipus, Hippolytus, and Agamemnon, but with little of the Greek tragic feeling for character and theme. The emphasis is on sensation and rhetoric, tending toward melodrama and bombast. The plays are of interest in this context mainly as the not entirely healthy inspiration for the precursors of Elizabethan tragedy in England.

The long hiatus in the history of tragedy between the Greeks and the Elizabethans has been variously explained. In the Golden Age of Roman literature, roughly from the birth of Virgil in 70 BC to the death of Ovid in AD 17, the Roman poets followed the example of Greek literature; although they produced great lyric and epic verse, their tragic drama lacked the probing freshness and directness fundamental to tragedy.

Reasons for the decline of Greek tragedy

Tragedy
and the
Christian
mass

With the collapse of the Roman world and the invasions of the barbarians came the beginnings of the long, slow development of the Christian Church. Churchmen and philosophers gradually forged a system, based on the Christian revelation, of the nature and destiny of man. The mass, with its daily reenactment of the sacrifice of Jesus Christ, its music, and its dramatic structure, may have provided something comparable to tragic drama in the lives of the people.

With the coming of the Renaissance, the visual arts more and more came to represent the afflictive aspects of life, and the word tragedy again came into currency. Chaucer (1340–1400) used the word in *Troilus and Criseyde*, and in *The Canterbury Tales* it is applied to a series of stories in the medieval style of *de casibus virorum illustrium*, meaning “the downfalls” (more or less inevitable) “of princes.” Chaucer used the word to signify little more than the turn of the wheel of fortune, against whose force no meaningful effort of man is possible. It remained for the Elizabethans to develop a theatre and a dramatic literature that reinstated the term on a level comparable to that of the Greeks.

Elizabethan. The long beginning of the Elizabethan popular theatre, like that of the Greek theatre, lay in religious ceremonies, probably in the drama in the liturgy of the two greatest events in the Christian year, Christmas and Easter. In the Early Church, exchanges between two groups of choristers, or between the choir and a solo voice, led to the idea of dialogue, just as it had in the development of Greek tragedy. The parts became increasingly elaborate, and costumes were introduced to individualize the characters. Dramatic gestures and actions were a natural development. More and more of the biblical stories were dramatized, much as the material of Homer was used by the Greek tragedians, although piously in this instance, with none of the tragic skepticism of the Greeks. In the course of generations, the popularity of the performances grew to such an extent that, to accommodate the crowds, they were moved, from inside the church to the porch, or square, in front of the church. The next step was the secularization of the management of the productions, as the towns and cities took them over. Day-long festivals were instituted, involving, as in the Greek theatre, the whole community. Cycles of plays were performed at York, Chester, and other English religious centres, depicting in sequences of short dramatic episodes the whole human story, from the Fall of Lucifer and the Creation to the Day of Doom. Each play was assigned to an appropriate trade guild (the story of Noah and the Ark, for example, went to the shipwrights), which took over complete responsibility for the production. Hundreds of actors and long preparation went into the festivals. These “miracle” and “mystery” plays, however crude they may now seem, dealt with the loftiest of subjects in simple but often powerful eloquence. Although the audience must have been a motley throng, it may well have been as involved and concerned as those of the Greek theatre.

Once the drama became a part of the secular life of the communities, popular tastes affected its religious orientation. Comic scenes, like those involving Noah’s nagging wife, a purely secular creation who does not appear in the Bible, became broader. The “tragic” scenes—anything involving the Devil or Doomsday—became more and more melodramatic. With the Renaissance came the rediscovery of the Greek and Roman cultures and the consequent development of a world view that led away from moral and spiritual absolutes and toward an increasingly skeptical individualism. The high poetic spirits of the mid-16th century began to turn the old medieval forms of the miracles and mysteries to new uses and to look to the ancient plays, particularly the lurid tragedies of Seneca, for their models. A bloody play, *Gorboduc*, by Thomas Sackville and Thomas Norton, first acted in 1561, is now known as the first formal tragedy in English, though it is far from fulfilling the high offices of the form in tone, characterization, and theme. Thomas Kyd’s *Spanish Tragedie* (c. 1589) continued the Senecan tradition of the “tragedy of blood” with somewhat more sophistication than *Gorboduc* but even more bloodletting. Elizabethan tragedy never

freed itself completely from certain melodramatic aspects of the influence of Seneca.

Marlowe and the first Christian tragedy. The first tragedian worthy of the tradition of the Greeks was Christopher Marlowe (1564–93). Of Marlowe’s tragedies, *Tamburlaine* (1587), *Doctor Faustus* (c. 1588), *The Jew of Malta* (1589), and *Edward II* (c. 1593), the first two are the most famous and most significant. In *Tamburlaine*, the material was highly melodramatic—Tamburlaine’s popular image was that of the most ruthless and bloody of conquerors. In a verse prologue, when Marlowe invites the audience to “View but his [Tamburlaine’s] picture in this tragic glass,” he had in mind little more, perhaps, than the trappings and tone of tragedy: “the stately tent of war,” which is to be his scene, and “the high astounding terms,” which will be his rhetoric. But he brought such imaginative vigour and sensitivity to bear that melodrama is transcended, in terms reminiscent of high tragedy. Tamburlaine, a Scythian shepherd of the 14th century, becomes the spokesman, curiously enough, for the new world of the Renaissance—iconoclastic, independent, stridently ambitious. Just as the Greek tragedians challenged tradition, Tamburlaine shouts defiance at all the norms, religious and moral, that Marlowe’s generation inherited. But Tamburlaine, although he is an iconoclast, is also a poet. No one before him on the English stage had talked with such magnificent lyric power as he does, whether it be on the glories of conquest or on the beauties of Zenocrate, his beloved. When, still unconquered by any enemy, he sickens and dies, he leaves the feeling that something great, however ruthless, has gone. Here once again is the ambiguity that was so much a part of the Greek tragic imagination—the combination of awe, pity, and fear that Aristotle defined.

In *Doctor Faustus* the sense of conflict between the tradition and the new Renaissance individualism is much greater. The claims of revealed Christianity are presented in the orthodox spirit of the morality and mystery plays, but Faustus’ yearnings for power over space and time are also presented with a sympathy that cannot be denied. Here is modern man, tragic modern man, torn between the faith of tradition and faith in himself. Faustus takes the risk in the end and is bundled off to hell in true mystery-play fashion. But the final scene does not convey that justice has been done, even though Faustus admits that his fate is just. Rather, the scene suggests that the transcendent human individual has been caught in the consequences of a dilemma that he might have avoided but that no imaginative man *could* have avoided. The sense of the interplay of fate and freedom is not unlike that of *Oedipus*. The sense of tragic ambiguity is more poignant in Faustus than in *Oedipus* or *Tamburlaine* because Faustus is far more introspective than either of the other heroes. The conflict is inner; the battle is for Faustus’ soul, a kind of conflict that neither the Greeks nor Tamburlaine had to contend with. For this reason, and not because it advocates Christian doctrine, the play has been called the first Christian tragedy.

Shakespearean tragedy. Shakespeare was a long time coming to his tragic phase, the six or seven years that produced his five greatest tragedies, *Hamlet* (c. 1601), *Othello* (c. 1602), *King Lear* (c. 1605), *Macbeth* (c. 1605), and *Antony and Cleopatra* (c. 1606). These were not the only plays written during those years. *Troilus and Cressida* may have come about the same time as *Hamlet*; *All’s Well That Ends Well*, shortly after *Othello*; and *Measure for Measure*, shortly before *King Lear*. But the concentration of tragedies is sufficient to distinguish this period from that of the comedies and history plays before and of the so-called romances afterward. Although the tragic period cannot entirely be accounted for in terms of biography, social history, or current stage fashions, all of which have been adduced as causes, certain questions should be answered, at least tentatively: What is Shakespeare’s major tragic theme and method? How do they relate to classical, medieval, and Renaissance traditions? In attempting to answer these questions, this proviso must be kept in mind: the degree to which he was consciously working in these traditions, consciously shaping his plays on early models, adapting Greek and Roman themes to his own purpose,

Conflict
between
tradition
and
individual-
ism

Cycles of
miracle
and mys-
tery plays

or following the precepts of Aristotle must always remain conjectural. On the one hand, there is the comment by Ben Jonson that Shakespeare had "small Latin and less Greek," and Milton in "L'Allegro" speaks of him as "fancy's child" warbling "his native wood-notes wild," as if he were unique, a sport of nature. On the other hand, Shakespeare knew Jonson (who knew a great deal of Latin and Greek) and is said to have acted in Jonson's *Sejanus* in 1603, a very classical play, published in 1605 with a learned essay on Aristotle as preface. It can be assumed that Shakespeare knew the tradition. Certainly the Elizabethan theatre could not have existed without the Greek and Roman prototype. For all of its mixed nature—with comic and melodramatic elements jostling the tragic—the Elizabethan theatre retained some of the high concern, the sense of involvement, and even the ceremonial atmosphere of the Greek theatre. When tragedies were performed, the stage was draped in black. Modern studies have shown that the Elizabethan theatre retained many ties with both the Middle Ages and the tradition of the Greeks.

Tragic elements in Shakespeare's comedies and histories

Shakespeare's earliest and most lighthearted plays reveal a sense of the individual, his innerness, his reality, his difference from every other individual, and, at times, his *plight*. Certain stock characters, to be sure, appear in the early comedies. Even Falstaff, that triumphant individual, has a prototype in the braggadocio of Roman comedy, and even Falstaff has his tragic side. As Shakespeare's art developed, his concern for the plight or predicament or dilemma seems to have grown. His earliest history plays, for instance (*Henry VI*, Parts I, II, III), are little more than chronicles of the great pageant figures—kingship in all its colour and potency. *Richard III*, which follows them, focusses with an intensity traditionally reserved for the tragic hero on one man and on the sinister forces, within and without, that bring him to destruction. From kingship, that is, Shakespeare turned to the king, the symbolic individual, the focal man, to whom whole societies look for their values and meanings. Thus Richard III is almost wholly sinister, though there exists a fascination about him, an all but tragic ambiguity.

Although Shakespeare's developing sense of the tragic cannot be summed up adequately in any formula, one might hazard the following: he progressed from the *individual* of the early comedies; to the *burdened* individual, such as, in *Henry IV*, Prince Hal, the future Henry V, who manipulates, rather than suffers, the tragic ambiguities of the world; and, finally, in the great tragedies, to (in one critic's phrase) the *overburdened* individual, Lear being generally regarded as the greatest example. In these last plays, man is at the limits of his sovereignty as a human being, where everything that he has lived by, stood for, or loved is put to the test. Like Prometheus on the crag, or Oedipus as he learns who he is, or Medea deserted by Jason, the Shakespearean tragic heroes are at the extremities of their natures. Hamlet and Macbeth are thrust to the very edge of sanity; Lear and, momentarily, Othello are thrust beyond it. In every case, as in the Greek plays, the destructive forces seem to combine inner inadequacies or evils, such as Lear's temper or Macbeth's ambition, with external pressures, such as Lear's "tiger daughters," the witches in *Macbeth*, or Lady Macbeth's importunity. Once the destructive course is set going, these forces operate with the relentlessness the Greeks called *Moirai*, or Fate.

Total vision of Shakespeare's tragedies

At the height of his powers, Shakespeare's tragic vision comprehended the totality of possibilities for good and evil as nearly as the human imagination ever has. His heroes are the vehicles of psychological, societal, and cosmic forces that tend to ennoble and glorify humanity or infect it and destroy it. The logic of tragedy that possessed him demanded an insistence upon the latter. Initially, his heroes make free choices and are free time after time to turn back, but they move toward their doom as relentlessly as did Oedipus. The total tragic statement, however, is not limited to the fate of the hero. He is but the centre of an action that takes place in a context involving many other characters, each contributing a point of view, a set of values or antivalues to the complex dialectic of the play. In Macbeth's demon-ridden Scotland, where weird things happen to men and horses turn cannibal, there is

the virtuous Malcolm, and society survives. Hamlet had the trustworthy friend Horatio, and, for all the blood-letting, what was "rotten" was purged. In the tragedies, most notably *Lear*, the Aeschylean notion of "knowledge through suffering" is powerfully dramatized; it is most obvious in the hero, but it is also shared by the society of which he is the focal figure. The flaw in the hero may be a *moral* failing or, sometimes, an excess of virtue; the flaw in society may be the rottenness of the Danish court in *Hamlet* or the corruption of the Roman world in *Antony and Cleopatra*; the flaw or fault or dislocation may be in the very universe itself, as dramatized by Lear's raving at the heavens or the ghosts that walk the plays or the witches that prophesy. All these faults, Shakespeare seems to be saying, are inevitabilities of the human condition. But they do not spell rejection, nihilism, or despair. The hero may die, but in the words of the novelist E.M. Forster to describe the redeeming power of tragedy, "he has given us life."

Such is the precarious balance a tragedian must maintain: the cold, clear vision that sees the evil but is not maddened by it, a sense of the good that is equally clear but refuses the blandishments of optimism or sentimentalism. Few have ever sustained the balance for long. Aeschylus tended to slide off to the right, Euripides to the left, and even Sophocles had his hero transfigured at Colonus. Marlowe's early death should perhaps spare him the criticism his first plays warrant. Shakespeare's last two tragedies, *Macbeth* and *Antony and Cleopatra*, are close to the edge of a valueless void. The atmosphere of *Macbeth* is murky with evil; the action moves with almost melodramatic speed from horror to horror. The forces for good rally at last, but Macbeth himself steadily deteriorates into the most nihilistic of all Shakespeare's tragic heroes, saved in nothing except the sense of a great nature, like Medea, gone wrong. *Antony*, in its ambiguities and irony, has been considered close to the Euripidean line of bitterness and detachment. Shakespeare himself soon modulated into another mood in his last plays, *Cymbeline* (c. 1609), *The Winter's Tale* (c. 1610), and *The Tempest* (c. 1611). Each is based on a situation that could have been developed into major tragedy had Shakespeare followed out its logic as he had done with earlier plays. For whatever reason, however, he chose not to. The great tragic questions are not pressed. *The Tempest*, especially, for all Prospero's charm and magnanimity, gives a sense of brooding melancholy over the ineradicable evil in mankind, a patient but sad acquiescence. All of these plays end in varying degrees of harmony and reconciliation. Shakespeare willed it so.

Decline in 17th-century England. From Shakespeare's tragedies to the closing of the theatres in England by the Puritans in 1642, the quality of tragedy is steadily worse, if the best of the Greek and Shakespearean tragedies are taken as a standard. Among the leading dramatists of the period—John Webster, Thomas Middleton, Francis Beaumont, John Fletcher, Cyril Tourneur, and John Ford—there were some excellent craftsmen and brilliant poets. Though each of them has a rightful place in the history of English drama, tragedy suffered a transmutation in their hands.

The Jacobean dramatists—those who flourished in England during the reign of James I—failed to transcend the negative tendencies they inherited from Elizabethan tragedy: a sense of defeat, a mood of spiritual despair implicit in Marlowe's tragic thought; in the nihilistic broodings of some of Shakespeare's characters in their worst moods—Hamlet, Gloucester in *Lear*, Macbeth; in the metaphoric implication of the theme of insanity, of man pressed beyond the limit of endurance, that runs through many of these tragedies; most importantly, perhaps, in the moral confusion ("fair is foul and foul is fair") that threatens to unbalance even the staunchest of Shakespeare's tragic heroes. This sinister tendency came to a climax about 1605 and was in part a consequence of the anxiety surrounding the death of Queen Elizabeth I and the accession of James I. Despite their negative tendencies, the Elizabethans, in general, had affirmed life and celebrated it; Shakespeare's moral balance, throughout even his darkest plays, remained firm. The Jacobean, on the

The Jacobean transmutation

other hand, were possessed by death. They became superb analysts of moral confusion and of the darkened vision of humanity at cross purposes, preying upon itself; of lust, hate, and intrigue engulfing what is left of beauty, love, and integrity. There is little that is redemptive or that suggests, as had Aeschylus, that evil might be resolved by the enlightenment gained from suffering. As in the tragedies of Euripides, the protagonist's margin of freedom grows ever smaller. "You are the deed's creature," cries a murderer to his unwitting lady accomplice in Middleton's *Changeling* (1622), and a prisoner of her deed she remains. Many of the plays maintained a pose of ironic, detached reportage, without the sense of sympathetic involvement that the greatest tragedians have conveyed from the beginning.

Some of the qualities of the highest tragedians have been claimed for John Webster. One critic points to his search for a moral order as a link to Shakespeare and sees in his moral vision a basis for renewal. Webster's *Duchess of Malfi* (c. 1613) has been interpreted as a final triumph of life over death. Overwhelmed by final unleashed terror, the Duchess affirms the essential dignity of man. Despite such vestiges of greatness, however, the trend of tragedy was downward. High moral sensitivity and steady conviction are required to resist the temptation to resolve the intolerable tensions of tragedy into either the comfort of optimism or the relaxed apathy of despair. Periods of the creation of high tragedy are therefore few and shortlived. The demands on artist and audience alike are very great. Forms wear out, and public taste seems destined to go through inevitable cycles of health and disease. What is to one generation powerful and persuasive rhetoric becomes bombast and bathos to the next. The inevitable materials of tragedy—violence, madness, hate, and lust—soon lose their symbolic role and become perverted to the uses of melodrama and sensationalism, mixed, for relief, with the broadest comedy or farce.

These corruptions had gone too far when John Milton, 29 years after the closing of the theatres, attempted to bring back the true spirit and tone of tragedy, which he called "the gravest, moralest, and most profitable of all other Poems." His *Samson Agonistes* (1671), however, is magnificent "closet tragedy"—drama more suitable for reading than for popular performance. Modelled on the *Prometheus*, it recalls Aeschylus' tragedy both in its form, in which the immobilized hero receives a sequence of visitors, and in its theme, in which there is a resurgence of the hero's spirit under stress. With Restoration comedy in full swing, however, and with the "heroic play" (an overly moralized version of tragedy) about to reach its crowning achievement in John Dryden's *All for Love* only seven years later (published 1678), *Samson Agonistes* was an anachronism.

Neoclassical. *Corneille and Racine.* Another attempt to bring back the ancient form had been going on for some time across the English Channel, in France. The French Classical tragedy, whose monuments are Pierre Corneille's *Cid* (1637) and Jean Racine's *Bérénice* (1670) and *Phèdre* (1677), made no attempt to be popular in the way of the Elizabethan theatre. The plays were written by and for intellectual aristocrats, who came together in an elite theatre, patronized by royalty and nobility. Gone were the bustle and pageantry of the Elizabethan tragedies, with their admixtures of whatever modes and moods the dramatists thought would work. The French playwrights submitted themselves to the severe discipline they derived from the Greek models and especially the "rules," as they interpreted them, laid down by Aristotle. The unities of place, time, and action were strictly observed. One theme, the conflict between Passion and Reason, was uppermost. The path of Reason was the path of Duty and Obligation (noblesse oblige), and that path had been clearly plotted by moralists and philosophers, both ancient and modern. In this sense there was nothing exploratory in the French tragedy; existing moral and spiritual norms were insisted upon. The norms are never criticized or tested as Aeschylus challenged the Olympians or as Marlowe presented, with startling sympathy, the Renaissance overreacher. Corneille's *Cid* shows Duty triumphant over Passion, and, as a reward, hero and heroine are happily united. By the

time of *Phèdre*, Corneille's proud affirmation of the power of the will and the reason over passion had given way to what Racine called "stately sorrow," with which he asks the audience to contemplate Phèdre's heroic, but losing, moral struggle. Her passion for her stepson, Hippolyte, bears her down relentlessly. Her fine principles and heroic will are of no avail. Both she and Hippolyte are destroyed. The action is limited to one terrible day; there is no change of scene; there is neither comic digression nor relief—the focus on the process by which a great nature goes down is sharp and intense. Such is the power of Racine's poetry (it is untranslatable), his conception of character, and his penetrating analysis of it, that it suggests the presence of Sophoclean "heroic humanism." In this sense it could be said that Racine tested the norms, that he uncovered a cruel injustice in the nature of a code that could destroy such a person as Phèdre. Once again, here is a world of tragic ambiguity, in which no precept or prescription can answer complicated human questions.

The English "heroic play." This ambiguity was all but eliminated in the "heroic play" that vied with the comedy of the Restoration stage in England in the latter part of the 17th century. After the vicissitudes of the Civil War, the age was hungry for heroism. An English philosopher of the time, Thomas Hobbes, defined the purpose of the type: "The work of an heroic poem is to raise admiration, principally for three virtues, valor, beauty, and love." Moral concern, beginning with Aeschylus, has always been central in tragedy, but in the works of the great tragedians this concern was exploratory and inductive. The moral concern of the heroic play is the reverse. It is deductive and dogmatic. The first rule, writes Dryden (following the contemporary French critic, René Le Bossu) in his preface to his *Troilus and Cressida* (1679), is "to make the moral of the work; that is, to lay down to yourself what that precept of morality shall be, which you would insinuate into the people . . ." In *All for Love* the moral is all too clear: Antony must choose between the path of honour and his illicit passion for Cleopatra. He chooses Cleopatra, and they are both destroyed. Only Dryden's poetry, with its air of emotional argumentation, manages to convey human complexities in spite of his moral bias and saves the play from artificiality—makes it, in fact, the finest near-tragic production of its age.

The eclipse of tragedy. Although the annals of the drama from Dryden onward are filled with plays called tragedies by their authors, the form as it has been defined here went into an eclipse during the late 17th, the 18th, and the early 19th centuries. Reasons that have been suggested for the decline include the politics of the Restoration in England; the rise of science and, with it, the optimism of the Enlightenment throughout Europe; the developing middle class economy; the trend toward reassuring deism in theology; and, in literature, the rise of the novel and the vogue of satire. The genius of the age was discursive and rationalistic. In France and later in England, belief in Evil was reduced to the perception of evils, which were looked upon as institutional and therefore remediable. The nature of man was no longer the problem; rather, it was the better organization and management of men. The old haunting fear and mystery, the sense of ambiguity at the centre of man's nature and of dark forces working against him in the universe, were replaced by a new and confident dogma. Tragedy never lost its high prestige in the minds of the leading spirits. Theorizing upon it were men of letters as diverse as Dr. Samuel Johnson, David Hume, Samuel Taylor Coleridge, and Percy Bysshe Shelley and German philosophers from Gotthold Lessing in the 18th century to Friedrich Nietzsche in the 19th. Revivals of Shakespeare's tragedies were often bowdlerized or altered, as in the happy ending for *Lear* in a production of 1681. Those who felt themselves called upon to write tragedies produced little but weak imitations. Shelley tried it once, in *The Cenci* (1819), but, as his wife wrote, "the bent of his mind went the other way"—which way may be seen in his *Prometheus Unbound* (1820), in which Zeus is overthrown and man enters upon a golden age, ruled by the power of love. Goethe had the sense to stay away from tragedy: "The mere attempt to write tragedy," he

Inductive
and
deductive
morality

The return
to Greek
models

said, "might be my undoing." He concluded his two-part *Faust* (1808, 1832) in the spirit of the 19th-century optimistic humanitarianism. It was not until the latter part of the 19th century, with the plays of a Norwegian, Henrik Ibsen, a Russian, Anton Chekhov, a Swede, August Strindberg, and, later, an American, Eugene O'Neill, that something of the original vision returned to inspire the tragic theatre.

A new vehicle: the novel. The theme and spirit of tragedy, meanwhile, found a new vehicle in the novel. This development is important, however far afield it may seem from the work of the formal dramatists. The English novelist Emily Brontë's *Wuthering Heights* (1847), in its grim Yorkshire setting, reflects the original concerns of tragedy: *i.e.*, the terrifying divisions in nature and human nature, love that creates and destroys, character at once fierce and pitiable, destructive actions that are willed yet seemingly destined, as if by a malicious fate, yet the whole controlled by an imagination that learns as it goes. Another English novelist, Thomas Hardy, in the preface to his *Woodlanders* (1887), speaks of the rural setting of this and other of his novels as being comparable to the stark and simple setting of the Greek theatre, giving his novels something of that drama's intensity and sharpness of focus. His grimly pessimistic view of man's nature and destiny and of the futility of human striving, as reflected in his novels *The Return of the Native* (1878), *Tess of the D'Urbervilles* (1891), and *Jude the Obscure* (1895), is barely redeemed for tragedy by his sense of the beauty of nature and of the beauty and dignity of human character and effort, however unavailing.

The work of the Polish-born English novelist Joseph Conrad (1857–1924) provides another kind of setting for novels used as vehicles of the tragic sense. *Lord Jim* (1900), originally conceived as a short story, grew to a full-length novel as Conrad found himself exploring in ever greater depth the perplexing, ambiguous problem of lost honour and guilt, expiation and heroism. Darkness and doubt brood over the tale, as they do over his long story "Heart of Darkness" (1899), in which Conrad's narrator, Marlow, again leads his listeners into the shadowy recesses of the human heart, with its forever unresolved and unpredictable capacities for good and evil.

Dostoyevsky's tragic view. In Russia, the novels of Fyodor Dostoyevsky, particularly *Crime and Punishment* (1866) and *The Brothers Karamazov* (1880), revealed a world of paradox, alienation, and loss of identity, prophetic of the major tragic themes of the 20th century. More than any earlier novelist, Dostoyevsky appropriated to his fictions the realm of the subconscious and explored in depth its shocking antinomies and discontinuities. Sigmund Freud, the founder of psychoanalysis, frequently acknowledged his indebtedness to Dostoyevsky's psychological insights. Dostoyevsky's protagonists are reminiscent of Marlow's Doctor Faustus, caught between the old world of orthodox belief and the new world of intense individualism, each with its insistent claims and justifications. The battleground is once more the soul of man, and the stakes are survival. Each of his major heroes—Raskolnikov in *Crime and Punishment* and the three Karamazovs—wins a victory, but it is in each case morally qualified, partial, or transient. The harmonious resolutions of the novels seem forced and are neither decisive of the action nor definitive of Dostoyevsky's total tragic view.

The American tragic novel. In America, Nathaniel Hawthorne's novel *The Scarlet Letter* (1850) and Herman Melville's *Moby Dick* (1851) are surprisingly complete embodiments of the tragic form, written as they were at a time of booming American optimism, materialistic expansion, and sentimentalism in fiction—and no tragic theatre whatever. In *The Scarlet Letter*, a story of adultery set in colonial New England, the heroine's sense of sin is incomplete; her spirited individualism insists (as she tells her lover) that "what we did had a consecration of its own." The resulting conflict in her heart and mind is never resolved, and, although it does not destroy her, she lives out her life in gray and tragic isolation. Melville said that he was encouraged by Hawthorne's exploration of "a certain tragic phase of humanity," by his deep broodings

and by the "blackness of darkness" in him, to proceed with similar explorations of his own in *Moby Dick*, which he dedicated to Hawthorne. Its protagonist, Captain Ahab, represents a return to what Melville called (defending Ahab's status as tragic hero) a "mighty pageant creature, formed for noble tragedies," whose "ponderous heart," "globular brain," and "nervous lofty language" prove that even an old Nantucket sea captain can take his place with kings and princes of the ancient drama. Shakespearean echoes abound in the novel; some of its chapters are written in dramatic form. Its theme and central figure, reminiscent of Job and Lear in their search for justice and of Oedipus in his search for the truth, all show what Melville might have been—a great tragic dramatist had there been a tragic theatre in America.

Some American novelists of the 20th century carried on, however partially, the tragic tradition. Theodore Dreiser's *American Tragedy* (1925) is typical of the naturalistic novel, which is also represented by the work of Stephen Crane, James T. Farrell, John Steinbeck. Though showing great sensitivity to environmental or sociological evils, such works fail to embody the high conception of character (as Melville describes it above) and are concerned mainly with externals, or reportage. The protagonists are generally "good" (or weak) and beaten down by society. The novels of Henry James, which span the period from 1876 to 1904, are concerned with what has been called the tragedy of manners. The society James projects is sophisticated, subtle, and sinister. The innocent and the good are destroyed, like Milly Theale in *The Wings of the Dove* (1902), who in the end "turns her face to the wall" and dies but in her death brings new vision and new values to those whose betrayals had driven her to her death.

The trend in American fiction, as in the drama, continued in the 20th century, toward the pathos of the victim—the somehow inadequate, the sometimes insignificant figure destroyed by such vastly unequal forces that the struggle is scarcely significant. F. Scott Fitzgerald's *Gatsby* in his novel *The Great Gatsby* (1925) is betrayed by his own meretricious dream, nurtured by a meretricious society. The hero of Ernest Hemingway's novel *A Farewell to Arms* (1929), disillusioned by war, makes a separate peace, deserts, and joins his beloved in neutral Switzerland. When she dies in childbirth, he sees it as still another example of how "they"—society, the politicians who run the war, or the mysterious forces destroying Catherine—get you in the end. The tone is lyric and pathetic rather than tragic (though Hemingway called the novel his *Romeo and Juliet*). Grief turns the hero away from, rather than toward, a deeper examination of life.

Only the novels of William Faulkner, in their range and depth and in their powerful assault on the basic tragic themes, recall unmistakably the values of the tragic tradition. His "saga of the South," as recounted in a series of novels (notably *Sartoris*, 1929; *The Sound and the Fury*, 1929; *As I Lay Dying*, 1930; *Sanctuary*, 1931; *Light in August*, 1932; *Absalom, Absalom!* 1936; *Intruder in the Dust*, 1948; *Requiem for a Nun*, 1951), incorporates some 300 years of Southern history from Indian days to the present. At first regarded as a mere exploiter of decadence, he can now be seen as gradually working beyond reportage and toward meaning. His sociology became more and more the "sin" of the South—the rape of the land, slavery, the catastrophe of the Civil War and its legacy of a cynical and devitalized materialism. Increasingly he saw the conflict as internal. The subject of art, Faulkner said in his 1949 Nobel Prize speech, is "the human heart in conflict with itself." His insistence is on guilt as the evidence of man's fate, and on the possibility of expiation as the assertion of man's freedom. Compassion, endurance, and the capacity to learn are seen to be increasingly effective in his characters. In the veiled analogies to Christ as outcast and redeemer in *Light in August* and in the more explicit Christology of *A Fable* (1954), in the pastoral serenity following the anguish and horror in *Light in August*, and in the high comedy of the last scene of *Intruder in the Dust*, Faulkner puts into tragic fiction the belief he stated in his Nobel speech: "I decline to accept the end of man."

Thomas
Hardy's
pessimism

Paradox,
alienation,
loss of
identity
in Dosto-
yevsky

Guilt and
expiation
in
Faulkner

TRAGEDY AND MODERN DRAMA

Tragic themes in Ibsen, Strindberg, and Chekhov. The movement toward naturalism in fiction in the latter decades of the 19th century did much to purge both the novel and the drama of the sentimentality and evasiveness that had so long emasculated them. In Norway Henrik Ibsen incorporated in his plays the smug and narrow ambitiousness of his society. The hypocrisy of overbearing men and women replace, in their fashion, the higher powers of the old tragedy. His major tragic theme is the futility, leading to catastrophe, of the idealist's effort to create a new and better social order. The "Problem play"—one devoted to a particular social issue—is saved in his hand from the flatness of a sociological treatise by a sense of doom, a pattern of retribution, reminiscent of the ancient Greeks. In *Pillars of Society* (1877), *The Wild Duck* (published 1884), *Rosmersholm* (published 1886), and *The Master Builder* (published 1892), for example, one sacrifice is expiated by another.

In Sweden, August Strindberg, influenced by Ibsen, was a powerful force in the movement. *The Father* (1887) and *Miss Julie* (published 1888) recall Ibsen's attacks on religious, moral, and political orthodoxies. Strindberg's main concern, however, is with the destructive effects of sexual maladjustment and psychic imbalance. Not since Euripides' *Medea* or Racine's *Phèdre* had the tragic aspects of sex come under such powerful analysis. In this respect, his plays look forward to O'Neill's.

Anton Chekhov, the most prominent Russian dramatist of the period, wrote plays about the humdrum life of inconspicuous, sensitive people (*Uncle Vanya*, 1899; *The Three Sisters*, 1901; and *The Cherry Orchard*, 1904, are typical), whose lives fall prey to the hollowness and tedium of a disintegrating social order. They are a brood of lesser Hamlets without his compensating vision of a potential greatness. As in the plays of the Scandinavian dramatists, Chekhov's vision of this social evil is penetrating and acute, but the powerful, resistant counterthrust that makes for tragedy is lacking. It is a world of victims.

American tragic dramatists. In little of the formal drama between the time of Ibsen, Strindberg, and Chekhov and the present are the full dimensions of tragedy presented. Some critics suggested that it was too late for tragedy, that modern man no longer valued himself highly enough, that too many sociological and ideological factors were working against the tragic temperament. The long and successful career of Eugene O'Neill may be a partial answer to this criticism. He has been called the first American to succeed in writing tragedy for the theatre, a fulfillment of his avowed purpose, for he had declared that in the tragic, alone, lay the meaning of life—and the hope. He sought in Freud's concept of the subconscious the equivalent of the Greek idea of fate and modelled his great trilogy, *Mourning Becomes Electra* (1931), on Aeschylus' *Oresteia*. Although the hovering sense of an ancient evil is powerful, the psychological conditioning controls the characters too nakedly. They themselves declare forces that determine their behaviour, so that they seem almost to connive in their own manipulation. *Desire Under the Elms* (1924) presents a harsh analysis of decadence in the sexual and avaricious intrigues of a New England farmer's family, unrelieved by manifestations of the transcendent human spirit. *The Great God Brown* (1926) and *Long Day's Journey into Night* (1939–41; first performance, 1956) come closer to true tragedy. In the latter, the capacity for self-knowledge is demonstrated by each member of the wrangling Tyrone family (actually, O'Neill's own; the play is frankly autobiographical). The insistent theme of the "death wish" (another example of Freud's influence), however, indicates too radical a pessimism for tragedy; even the character of Edmund Tyrone, O'Neill's own counterpart, confesses that he has always been a little in love with death, and in another late play, *The Iceman Cometh* (1939), the death wish is more strongly expressed. Although he never succeeded in establishing a tragic theatre comparable to the great theatres of the past, O'Neill made a significant contribution in his sustained concentration on subjects at least worthy of such a theatre. He made possible the significant, if slighter, contributions of

O'Neill:
pessimism
too radical
for tragedy

Arthur Miller, whose *Death of a Salesman* (1949) and *A View from the Bridge* (1955) contain material of tragic potential that is not fully realized. Tennessee Williams' *Streetcar Named Desire* (1947) is a sensitive study of the breakdown of a character under social and psychological stress. As with Miller's plays, however, it remains in the area of pathos rather than tragedy.

Other serious drama. The 20th century has produced much serious and excellent drama, which, though not in the main line of the tragic tradition, deserves mention. In British theatre, George Bernard Shaw's *Saint Joan* (1923) and T.S. Eliot's *Murder in the Cathedral* (1935) dramatized with great power both doubt and affirmation, the ambiguity of human motives, and the possibility of fruitless suffering that are true of the human condition as reflected by tragedy. During the Irish literary revival, the work of J.M. Synge (*Riders to the Sea*, 1904) and Sean O'Casey (*Shadow of a Gunman*, 1923), like Faulkner's work, sought a tragic theme in the destiny of a whole people. The masterpiece of this movement, however, is not a tragedy but a comic inversion of the ancient tragedy of *Oedipus*—Synge's *Playboy of the Western World* (1907).

The drama of social protest—exemplified in such works as the Russian Maksim Gorky's *Lower Depths* (1902), the German Bertolt Brecht's *Threepenny Opera* (1928) and *Mother Courage* (1941), and the American Clifford Odets' *Waiting for Lefty* (1935)—shares the tragedians' concern for evils that frustrate or destroy human values. The evils, however, are largely external, identifiable, and, with certain recommended changes in the social order, remediable. The type shows how vulnerable tragedy is to dogma or programs of any sort. A British author, George Orwell, suggested in *Nineteen Eighty-four* that tragedy would cease to exist under pure Marxist statism. Brecht's fine sense of irony and moral paradox redeem him from absolute dogmatism but give his work a hard satire thrust that is inimical to tragedy. Traditional values and moral imperatives are all but neutralized in the existentialist worlds of the dramas and novels of Jean-Paul Sartre and Albert Camus, two outstanding philosopher-dramatists of the post World War II era. In their works, the protagonist is called upon to forge his own values, if he can, in a world in which the disparity between the ideal (what man longs for) and the real (what he gets) is so great as to reduce the human condition to incoherence and absurdity. Plays that led to the coinage of the term the theatre of the absurd are exemplified by *Waiting for Godot* (1952) and *The Killer* (1959), respectively by the Irish writer Samuel Beckett and the Romanian Eugène Ionesco, both of whom pursued their careers in Paris. Here, the theme of victimization is at its extreme, the despair and defeat almost absolute.

A coherent and affirmative view of man, society, and the cosmos is vital to tragedy—however tentative the affirmation may be. Unresolved questions remain at the end of every tragedy. There is always an irrational factor, disturbing, foreboding, not to be resolved by the sometime consolations of philosophy and religion or by any science of the mind or body; there is irretrievable loss, usually though not necessarily symbolized by the death of the hero. In the course of the action, however, in the development of character, theme, and situation and in the conceptual suggestiveness of language, tragedy presents the positive terms in which these questions might be answered. The human qualities are manifest, however limited; man's freedom is real, however marginal. The forces that bear him down may be mysterious but actual—fate, the gods, chance, the power of his own or the race's past working through his soul. Though never mastered, they can be contended with, defied, and, at least in spirit, transcended. The process is cognitive; man can learn.

Absence of tragedy in Oriental drama. In no way can the importance of a conceptual basis for tragedy be better illustrated than by a look at other drama-producing cultures with radically different ideas of the individual, his nature, and his destiny. While the cultures of India, China, and Japan have produced significant and highly artistic drama, there is little here to compare in magnitude, intensity, and freedom of form to the tragedies of the West.

In Buddhist teaching, the aim of the individual is to sup-

Tragedy
in works
of social
protest

press and regulate all those questioning, recalcitrant, rebellious impulses that first impel the Western hero toward his tragic course. The goal of Nirvāṇa is the extinction of those impulses, the quieting of the passions, a kind of *quietus* in which worldly existence ceases. Western tragedy celebrates life, and the tragic hero clings to it: to him, it is never "sweet to die" for his country or for anything else, and the fascination for Western audiences is to follow the hero—as it were, *from the inside*—as he struggles to assert himself and his values against whatever would deny them. In Oriental drama, there is no such intense focus on the individual. In the Japanese Nō plays, for instance, the hero may be seen in moments of weariness and despair, of anger or confusion, but the mood is lyric, and the structure of the plays is ritualistic, with a great deal of choral intoning, dancing, and stylized action. Although a number of Nō plays can be produced together to fill a day's performance, the individual plays are very short, hardly the length of a Western one-act play. Nō plays affirm orthodoxy, rather than probing and questioning it, as Western tragedies do.

The drama in India has a long history, but there too the individual is subordinated to the mood of the idyll or romance or epic adventure. Perhaps one reason why the drama of India never developed the tragic orientation of the West is its removal from the people; it has never known the communal involvement of the Greek and Elizabethan theatres. Produced mainly for court audiences, an upper class elite, it never reflected the sufferings of common (or uncommon) humanity. Only recently has the drama in China embraced the vigour and realism of the common people, but the drama is in the service not of the individual but of a political ideology, which replaces the traditional themes of ancestor worship and filial piety. In all this, the mighty pageant figure—Oedipus, Prometheus, Lear, or Ahab standing for the individual as he alone sees and feels the workings of an unjust universe—is absent.

Tragedy in
Nō drama

An example from the Nō plays will illustrate these generalizations. In *The Hoka Priests*, by Zenchiku Ujino (1414–99), a son is confronted with Hamlet's problem—*i.e.*, that of avenging the death of his father. He is uncertain how to proceed, since his father's murderer has many bold fellows to stand by him, while he is all alone. He persuades his brother, a priest, to help him, and disguising themselves as priests, they concoct a little plot to engage the murderer in religious conversation. There are a few words of lament—"Oh why./ Why back to the bitter World/ Are we borne by our intent?"—and the Chorus sings lyrically about the uncertainties of life. The theme of the conversation is the unreality of the World and the reality of Thought. At an appropriate moment, the brothers cry, "Enough! Why longer hide our plot?" The murderer places his hat on the floor and exits. The brothers mime the killing of the murderer in a stylized attack upon the hat, while the Chorus describes and comments on the action: "So when the hour was come/ Did these two brothers/ By sudden resolution/ Destroy their father's foe./ For valour and piety are their names remembered/ Even in this aftertime" (translated by Arthur Waley, *The Nō Plays of Japan*, 1921).

Thus the Nō avoids directly involving the audience in the emotions implicit in the events portrayed on the stage. It gives only a slight hint of the spiritual struggle in the heart of the protagonist—a struggle that is always speedily resolved in favour of traditional teaching. In play after play the action does not take place before our eyes but is reenacted by the ghost of one of the participants. Thus, the events presented are tinged with memory or longing—hardly the primary emotions that surge through and invigorate Western tragedy at its best.

Loss of viability in the West. The absence, even in the West, of a great tragic theatre in the 20th century may be explained by the pantheon of panaceas to which modern man has subscribed. Politics, psychology, social sciences, physical sciences, nationalism, the occult—each offered a context in terms of which he might act out his destiny, were it not crowded out by the others. Modern man is not tested but harried and not by gods but, too often, by demons. In the dramas of Athens and England, tragedy

was born of the impossibility of a clear-cut victory in man's struggle with powers greater than himself. In the modern drama, the struggle itself seems impossible.

The would-be hero is saved from a meaningful death by being condemned to a meaningless life. This, too, however, has its tragic dimension, in its illustration of the power of evil to survive from millennium to millennium in the presence or the absence of the gods.

Tragedy is a means of coming to terms with that evil. To assume that tragedy has lost viability is to forget that this viability was seriously questioned by the first Western philosopher to address himself to the problem. An account of the development of the theory of tragedy will reveal a resourcefulness in man's critical powers that can help to compensate, or occasionally even supersede, his lapsing creative powers. (R.B.S.)

THEORY OF TRAGEDY

Classical theories. As the great period of Athenian drama drew to an end at the beginning of the 4th century BC, Athenian philosophers began to analyze its content and formulate its structure. In the thought of Plato (c. 427–347 BC), the history of the criticism of tragedy began with speculation on the role of censorship. To Plato (in the dialogue on the *Laws*) the state was the noblest work of art, a representation (*mimēsis*) of the fairest and best life. He feared the tragedians' command of the expressive resources of language, which might be used to the detriment of worthwhile institutions. He feared, too, the emotive effect of poetry, the Dionysian element that is at the very basis of tragedy. Therefore, he recommended that the tragedians submit their works to the rulers, for approval, without which they could not be performed. It is clear that tragedy, by nature exploratory, critical, independent, could not live under such a regimen.

Plato is answered, in effect and perhaps intentionally, by Aristotle's *Poetics*. Aristotle (384–322 BC) defends the purgative power of tragedy and, in direct contradiction to Plato, makes moral ambiguity the essence of tragedy. The tragic hero must be neither a villain nor a virtuous man but a "character between these two extremes, . . . a man who is not eminently good and just, yet whose misfortune is brought about not by vice or depravity, but by some error or frailty [*hamartia*]." The effect on the audience will be similarly ambiguous. A perfect tragedy, he says, should imitate actions that excite "pity and fear." He uses Sophocles' *Oedipus the King* as a paradigm. Near the beginning of the play, Oedipus asks how his stricken city (the counterpart of Plato's state) may cleanse itself, and the world he uses for the purifying action is a form of the word catharsis. The concept of catharsis provides Aristotle with his reconciliation with Plato, a means by which to satisfy the claims of both ethics and art. "Tragedy," says Aristotle, "is an imitation [*mimēsis*] of an action that is serious, complete and of a certain magnitude . . . through pity and fear effecting the proper purgation [catharsis] of these emotions." Ambiguous means may be employed. Aristotle maintains in contrast to Plato, to a virtuous and purifying end.

To establish the basis for a reconciliation between ethical and artistic demands, Aristotle insists that the principal element in the structure of tragedy is not character but plot. Since the erring protagonist is always in at least partial opposition to the state, the importance of tragedy lies not in him but in the enlightening event. "Most important of all," Aristotle said, "is the structure of the incidents. For tragedy is an imitation not of men but of an action and of life, and life consists in action, and its end is a mode of action, not a quality . . ." Aristotle considered the plot to be the soul of a tragedy, with character in second place. The goal of tragedy is not suffering but the knowledge that issues from it, as the denouement issues from a plot. The most powerful elements of emotional interest in tragedy, according to Aristotle, are reversal of intention or situation (*peripeteia*) and recognition scenes (*anagnōrīsis*), and each is most effective when it is coincident with the other. In *Oedipus*, for example, the messenger who brings Oedipus news of his real parentage, intending to allay his fears, brings about a sudden reversal of his fortune, from

Ethics
and art—
Plato and
Aristotle

happiness to misery, by compelling him to recognize that his wife is also his mother.

Later critics found justification for their own predilections in the authority of Greek drama and Aristotle. For example, the Roman poet Horace (65–8 BC), in his *Ars poetica* (*Art of Poetry*), elaborated the Greek tradition of extensively narrating offstage events into a dictum on decorum forbidding events such as Medea's butchering of her boys from being performed on stage. And where Aristotle had discussed tragedy as a separate genre, superior to epic poetry, Horace discussed it as a genre with a separate style, again with considerations of decorum foremost. A theme for comedy may not be set forth in verses of tragedy; each style must keep to the place allotted it.

On the basis of this kind of stylistic distinction, the *Aeneid*, the epic poem of Virgil, Horace's contemporary, is called a tragedy by the fictional Virgil in Dante's *Divine Comedy*, on the grounds that the *Aeneid* treats only of lofty things. Dante (1265–1321) calls his own poem a comedy partly because he includes "low" subjects in it. He makes this distinction in his *De vulgari eloquentia* (1304–05; "Of Eloquence in the Vulgar") in which he also declares the subjects fit for the high, tragic style to be salvation, love, and virtue. Despite the presence of these subjects in this poem, he calls it a comedy because his style of language is "careless and humble" and because it is in the vernacular tongue rather than Latin. Dante makes a further distinction:

Comedy . . . differs from tragedy in its subject matter, in this way, that tragedy in its beginning is admirable and quiet, in its ending or catastrophe foul and horrible. . . . From this it is evident why the present work is called a comedy.

Dante's emphasis on the outcome of the struggle rather than on the nature of the struggle is repeated by Chaucer and for the same reason: their belief in the providential nature of human destiny. Like Dante, he was under the influence of *De consolazione philosophiae* (*Consolation of Philosophy*), the work of the 6th-century Roman philosopher Boethius (c. 480–524) that he translated into English. Chaucer considered Fortune to be beyond the influence of the human will. In his *Canterbury Tales*, he introduces "The Monk's Tale" by defining tragedy as "a certeyn storie . . . / of him that stood in greet prosperitee, / And is y-fallen out of heigh degree / Into miserie, and endeth wrecchedly." Again, he calls his *Troilus and Criseyde* a tragedy because, in the words of Troilus, "all that comth, comth by necessitee . . . / That forsyght of divine purveyaunce / Hath seyn alwey me to forgon Criseyde."

Elizabethan approaches. The critical tradition of separating the tragic and comic styles is continued by the Elizabethan English poet Sir Philip Sidney, whose *Defence of Poesie* (also published as *An Apologie for Poetrie*) has the distinction of containing the most extended statement on tragedy in the English Renaissance and the misfortune of having been written in the early 1580s (published 1595), before the first plays of Shakespeare, or even of Marlowe. Nevertheless, Sidney wrote eloquently of "high and excellent tragedy, that . . . with stirring the affects of admiration and commiseration teacheth the uncertainty of this world and upon how weak foundations gilden roofs are builded."

Since the word admiration here means awe, Sidney's "admiration and commiseration" are similar to Aristotle's "pity and fear." He differs from Aristotle, however, in preferring epic to tragic poetry. The Renaissance was almost as concerned as Plato with the need to justify poetry on ethical grounds, and Sidney ranks epic higher than tragedy because it provides morally superior models of behaviour.

Sidney goes further than mere agreement with Aristotle, however, in championing the unities of time and place. Aristotle had asserted the need for a unity of time: "Tragedy endeavors, as far as possible, to confine itself to a single revolution of the sun, or but slightly to exceed this limit." Sidney, following the lead of a 16th-century Italian Neoclassicist, Ludovico Castelvetro, added the unity of place: "the stage should always represent but one place, and the uttermost time presupposed in it should be, both by Aristotle's precept and common reason, but one day . . ." Sidney also seconds Horace's disapproval of the

mingling of styles, which Sidney says produces a "mongrel tragicomedy."

Shakespeare's opinion of the relative merits of the genres is unknown, but his opinion of the problem itself may be surmised. In *Hamlet* he puts these words in the mouth of the foolish old pedant Polonius: "The best actors in the world, either for tragedy, comedy, history, pastoral, pastoral-comical, historical-pastoral, tragical-historical, tragical-comical-historical-pastoral: scene indivisible, or poem unlimited . . ." (Act II, scene 2). As to the classical unities, Shakespeare adheres to them only twice and neither time in a tragedy, in *The Comedy of Errors* and *The Tempest*. And through the mouths of his characters, Shakespeare, like Aristotle, puts himself on both sides of the central question of tragic destiny—that of freedom and necessity. Aristotle says that a tragic destiny is precipitated by the hero's tragic fault, his "error or frailty" (hamartia), but Aristotle also calls this turn of events a change of "fortune." Shakespeare's Cassius in *Julius Caesar* says, "The fault, dear Brutus, is not in our stars, / But in ourselves . . .," and in *King Lear*, Edmund ridicules a belief in fortune as the "foppery of the world." But Hamlet, in a comment on the nature of hamartia, is a fatalist when he broods on the "mole of nature," the "one defect" that some men are born with, "wherein they are not guilty," and that brings them to disaster (Act I, scene 4). Similarly, Sophocles' Oedipus, though he says, "It was Apollo who brought my woes to pass," immediately adds, "it was my hand that struck my eyes." These ambiguities are a powerful source of the tragic emotion of Athenian and Elizabethan drama, unequalled by traditions that are more sure of themselves, such as French Neoclassicism, or less sure of themselves, such as 20th-century drama.

Neoclassical theory. In the Neoclassical period Aristotle's reasonableness was replaced by rationality, and his moral ambiguity by the mechanics of "poetic justice." In the 17th century, under the guise of a strict adherence to Classical formulas, additional influences were brought to bear on the theory of tragedy. In France, the theological doctrine of Jansenism, which called for an extreme orthodoxy, exercised a strong influence. In England, the restoration of the monarchy in 1660, with the reopening of the theatres, introduced a period of witty and lusty literature. In both nations, the influence of natural law—the idea that laws binding upon humanity are inferable from nature—increased, along with the influence of the exact sciences. Critics in both nations declared that Aristotle's "rules" were made to reduce nature into a method.

In his 1679 preface to Shakespeare's *Troilus and Cressida*, Dryden says, "we lament not, but detest a wicked man, we are glad when we behold his crimes are punished, and that Poetical justice is done upon him." Similar sentiments, calling for the punishment of crimes and the reward of virtue, were expressed in France. Catharsis had become vindication. Thomas Rymer, one of the most influential English critics of the time, in *The Tragedies of The Last Age* (1678), wrote that

besides the purging of the passions, something must stick by observing . . . that necessary relation and chain, whereby the causes and the effects, the virtues and rewards, the vices and their punishments are proportion'd and link'd together, how deep and dark soever are laid the Springs, and however intricate and involv'd are their operations.

The effect was to rob tragedy of a great deal of its darkness and depth. The temper of the age demanded that mystery be brought to the surface and to the light, a process that had effects not merely different from but in part antipathetic to tragedy. Nicolas Boileau, the chief spokesman of the French Neoclassical movement, in his discussion of pity and fear in *Art Poétique* (1674), qualified these terms with the adjectives "beguiling" and "pleasant" (*pitié charmante, douce terreur*), which radically changed their meaning. The purged spectator became a grateful patient. In his preface to *Phèdre* (1677), Racine subscribed to the *quid pro quo* view of retribution.

I have written no play in which virtue has been more celebrated than in this one. The smallest faults are here severely punished; the mere idea of a crime is looked upon with as much horror as the crime itself.

Fate and will in comedy and tragedy

Vindication rather than catharsis

Of Phèdre herself, his greatest heroine, he says,

I have taken the trouble to make her a little less hateful than she is in the ancient versions of this tragedy, in which she herself resolves to accuse Hippolytus. I judged that that calumny had about it something too base and black to be put into the mouth of a Princess. . . . This depravity seemed to me more appropriate to the character of a nurse, whose inclinations might be supposed to be more servile. . . .

For Aristotle, pity and fear made a counterpoint typical of Classicism, each tempering the other to create a balance. For Racine, pity and fear each must be tempered in itself. In the marginalia to his fragmentary translation of Aristotle's *Poetics*, Racine wrote that in arousing the passions of pity and fear, tragedy

removes from them whatever they have of the excessive and the vicious and brings them back to a moderated condition and conformable to reason.

Corneille contradicted Aristotle outright. Discussing *Le Cid* he said, in *A Discourse on Tragedy* (1660),

Our pity ought to give us fear of falling into similar misfortune, and purge us of that excess of love which is the cause of their disaster. . . . but I do not know that it gives us that, or purges us, and I am afraid that the reasoning of Aristotle on this point is but a pretty idea. . . . it is not requisite that these two passions always serve together. . . . it suffices. . . . that one of the two bring about the purgation. . . .

The accommodation of tragedy to Neoclassical ideas of order demanded a simplification of tragedy's complexities and ambiguities. The simplifying process was now inspired, however, by the fundamental tenet of all primitive scientific thought namely, that orderliness and naturalness are in a directly proportionate relationship. Racine declared the basis of the naturalistic effect in drama to be a strict adherence to the unities, which now seem the opposite of naturalistic. In his preface to *Bérénice* (1670), he asked what probability there could be when a multitude of things that would scarcely happen in several weeks are made to happen in a day. The illusion of probability, which is the Aristotelian criterion for the verisimilitude of a stage occurrence, is made to sound as if it were the result of a strict dramaturgical determinism, on the grounds that necessity is the truest path to freedom.

Racine and Corneille both contradicted Dante and Chaucer on the indispensability of a catastrophic final scene. "Blood and deaths," said Racine, are not necessary, for "it is enough that the action be grand, that the actors be heroic, that the passions be aroused" to produce "that stately sorrow that makes the whole pleasure of tragedy" (preface to *Bérénice*).

Milton was artistically much more conservative. He prefaced his *Samson Agonistes* (1671) with a warning against the

error of intermixing Comic stuff with Tragic sadness and gravity; or introducing trivial and vulgar persons: which by all judicious hath been counted absurd; and brought in without discretion, corruptly to gratify the people.

He bypassed Shakespeare for the ancients and ranked Aeschylus, Sophocles, and Euripides as tragic poets unequalled yet by any others. Part of the rule, for Milton, was that which affirmed the unities. In his concurrence with the Classical idea of the purgative effect of pity and fear, Milton combined reactionary aesthetics with the scientific spirit of the recently formed Royal Society.

Nor is Nature wanting in her own effects to make good his assertion [Aristotle on catharsis]: for so, in Phisic things of melancholic hue and quality are used against melancholy, sour against sour, salt to remove salt humours.

Dryden spoke against a delimiting conception of either the genres or the unities. Speaking in the guise of Neander in *Of Dramatick Poesie, an Essay* (1668), he said that it was

to the honour of our nation, that we have invented, increased, and perfected a more pleasant way of writing for the stage, than was ever known to the ancients or moderns of any nation, which is tragi-comedy.

The French dramatists, he felt, through their observance of the unities of time and place, wrote plays characterized by a dearth of plot and narrowness of imagination. Ra-

cine's approach to the question of probability was turned completely around by Dryden, who asked:

How Many beautiful accidents might naturally happen in two or three days, which cannot arrive with any probability in the compass of twenty-four hours?

The definitive critique of Neoclassical restrictions was not formulated, however, until the following century, when it was made by Samuel Johnson and was, significantly, part of his 1765 preface to Shakespeare, the first major step in the long process of establishing Shakespeare as the preeminent tragic poet of post-Classical drama. On genre he wrote:

Shakespeare's plays are not in the rigorous and critical sense either tragedies or comedies, but compositions of a distinct kind: . . . expressing the course of the world, in which the loss of one is the gain of another; in which, at the same time, the reveller is hasting to his wine, and the mourner burying his friend. . . . That this is a practice contrary to the rules of criticism will be readily allowed; but there is always an appeal open from criticism to nature.

And on the unities:

The necessity of observing the unities of time and place arises from the supposed necessity of making the drama credible. [But] the objection arising from the impossibility of passing the first hour at Alexandria, and the next at Rome, supposes, that when the play opens, the spectator really imagines himself at Alexandria. . . . Surely he that imagines this may imagine more.

Johnson's appeal to nature was the essence of subsequent Romantic criticism.

Romantic theories. Lessing was the first important Romantic critic. He stated one of Romanticism's chief innovations in his *Hamburg Dramaturgy* (1767-69):

The names of princes and heroes can lend pomp and majesty to a play, but they contribute nothing to our emotion. The misfortune of those whose circumstances most resemble our own, must naturally penetrate most deeply into our hearts, and if we pity kings, we pity them as human beings, not as kings.

Within a generation, revolutions in Europe and America offered social expression of this literary precept, and a dramatic tradition dominant for 22 centuries was upturned. From the time of Aristotle, who thought that the tragic hero should be highly renowned and prosperous, the tragic hero had been an aristocrat, if not a man of royal blood. With the exception of their minor or peripheral characters, the tragic dramas of Athens, England, and France told nothing of the destinies of the mass of mankind. All this was now changed.

But it is not certain that what was good for the revolution was good for tragedy. Coleridge in his critical writings of 1808-18 said that:

there are two forms of disease most preclusive of tragic worth. The first [is] a sense and love of the ludicrous, and a diseased sensibility of the assimilating power. . . . that in the boldest bursts of passion will lie in wait, or at once kindle into jest. . . . The second cause is matter of exultation to the philanthropist and philosopher, and of regret to the poet. . . . namely, the security, comparative equality, and ever-increasing sameness of human life.

In accord with this distaste for an excess of the mundane, Coleridge attacked the new German tragedies in which "the dramatist becomes a novelist in his directions to the actors, and degrades tragedy to pantomime." To describe, or rather indicate, what tragedy should ideally be, Coleridge said "it is not a copy of nature; but it is an imitation."

Coleridge's operative words and phrases in his discussions of tragedy were "innate," "from within," "implicit," "the being within," "the inmost heart," "our inward nature," "internal emotions," and "retired recesses." The new philosophical dispensation in Coleridge, like the new social dispensation in Lessing, reversed the old priorities; and where there were once princes there were now burghers, and where there were once the ordinances of God and the state there were now the dictates of the heart. By means of this reversal, Coleridge effected a reconciliation of the "tragedy of fate" and the "tragedy of character" in his description of the force of fate as merely the embodiment

Johnson's
definitive
critique

Coleridge's
reconciliation
of
fate and
character

of an interior compulsion different in scale but not in kind from the interior compulsions of character. In Classical tragedy, he said the human "will" was "exhibited as struggling with fate, a great and beautiful instance and illustration of which is the Prometheus of Aeschylus; and the deepest effect is produced, when the fate is represented as a higher and intelligent will. . . ."

According to Coleridge, Shakespeare used the imaginative "variety" that characterizes man's inward nature in place of the mechanical regularity of the Neoclassical unities to produce plays that were "neither tragedies nor comedies, nor both in one, but a different genus, diverse in kind, not merely different in degree,—romantic dramas or dramatic romances." In his preoccupation with the mixture of genres and his distinction between the "mechanical" (Neoclassicism) and the "organic" (Shakespeare), Coleridge was influenced by *Lectures on Dramatic Art and Literature* (delivered 1808–09, published 1809–11), by August Wilhelm von Schlegel, perhaps the most influential of German Romantic critics.

Like Coleridge and most Romantic critics of tragedy, Schlegel found his champion in Shakespeare, and, also like them, he was preoccupied with the contrast between Classic and Romantic. Like Coleridge, Schlegel emphasized Shakespeare's inwardness, what Coleridge called his "implicit wisdom deeper even than our consciousness." It is in Shakespeare's most profound insights that Schlegel locates one of the principal distinctions between Classical and Shakespearean tragedy, in what he calls Shakespeare's "secret irony." The irony in *Oedipus the King* consists in the relation between the audience's knowledge of the protagonist's situation and his own ignorance of it. But Shakespeare's "readiness to remark the mind's fainter and involuntary utterances" is so great, says Schlegel, that "nobody ever painted so truthfully as he has done the facility of self-deception, the half self-conscious hypocrisy towards ourselves, with which even noble minds attempt to disguise the almost inevitable influence of selfish motives in human nature."

The irony Schlegel sees in Shakespeare's characterizations also extends to the whole of the action, as well as to the separate characters. In his discussion of it he suggests the reason for the difficulty of Shakespeare's plays and for the quarrelsome, irreconcilable "interpretations" among Shakespeare's commentators:

Most poets who portray human events in a narrative or dramatic form take themselves apart, and exact from their readers a blind approbation or condemnation of whatever side they choose to support or oppose. . . . When, however, by a dexterous manoeuvre, the poet allows us an occasional glance at the less brilliant reverse of the medal, then he makes, as it were, a sort of secret understanding with the select circle of the more intelligent of his readers or spectators; he shows them that he had previously seen and admitted the validity of their tacit objections; that he himself is not tied down to the represented subject but soars freely above it. . . .

In Greek tragedy, the commentary by the chorus was an explicit and objective fact of the drama itself. In the presentation of Shakespeare's plays, such a commentary is carried on in the separate minds of the spectators, where it is diffused, silent, and not entirely sure of itself. When the spectators speak their minds after the curtain falls, it is not surprising that they often disagree.

In *Oedipus the King*, which Aristotle cited as the model of Classical tragedy, the irony of the protagonist's situation is evident to the spectator. In *Hamlet*, however, according to the American philosopher George Santayana, writing in 1908, it is the secret ironies, half-lights, and self-contradictions that make it the central creation of Romantic tragedy. As has been noted, Coleridge objected to the dramatist's giving directions to the actors, but part of the price of not having them is to deny to the audience as well an explicit indication of the playwright's meaning.

George Wilhelm Friedrich Hegel (1770–1831), the immensely influential German philosopher, in his *Aesthetik* (1820–29), proposed that the sufferings of the tragic hero are merely a means of reconciling opposing moral claims. The operation is a success because of, not in spite of, the fact that the patient dies. According to Hegel's account of

Greek tragedy, the conflict is not between good and evil but between goods that are each making too exclusive a claim. The heroes of ancient tragedy, by adhering to the *one* ethical system by which they molded their own personality, must come into conflict with the ethical claims of another. It is the moral one-sidedness of the tragic actor, not any negatively tragic fault in his morality or in the forces opposed to him, that proves his undoing, for both sides of the contradiction, if taken by themselves, are justified.

The nuclear Greek tragedy for Hegel is, understandably, Sophocles' *Antigone*, with its conflict between the valid claims of conscience (Antigone's obligation to give her brother a suitable burial) and law (King Creon's edict that enemies of the state should not be allowed burial). The two claims represent what Hegel regards as essentially concordant ethical claims. Antigone and Creon are, in this view, rather like pawns in the Hegelian dialectic—his theory that thought progresses from a thesis (*i.e.*, an idea), through an antithesis (an idea opposing the original thesis), to a synthesis (a more comprehensive idea that embraces both the thesis and antithesis), which in turn becomes the thesis in a further progression. At the end of *Antigone*, something of the sense of mutually appeased, if not concordant, forces does obtain after Antigone's suicide and the destruction of Creon's family. Thus, in contrast to Aristotle's statement that the tragic actors should represent not an extreme of good or evil but something between, Hegel would have them too good to live; that is, too extreme an embodiment of a particular good to survive in the world. He also tends to dismiss other traditional categories of tragic theory. For instance, he prefers his own kind of catharsis to Aristotle's—the feeling of *reconciliation*.

Hegel's emphasis on the correction of moral imbalances in tragedy is reminiscent of the "poetic justice" of Neoclassical theory, with its similar dialectic of crime and punishment. He sounds remarkably like Racine when he claims that, in the tragic denouement, the necessity of all that has been experienced by particular individuals is seen to be in complete accord with reason and is harmonized on a true ethical basis. But where the Neoclassicists were preoccupied with the unities of time and place, Hegel's concerns, like those of other Romantics, are inward. For him, the final issue of tragedy is not the misfortune and suffering of the tragic antagonists but rather the satisfaction of spirit arising from "reconciliation." Thus, the workings of the spirit, in Hegel's view, are subject to the rationalistic universal laws.

Hegel's system is not applicable to Shakespearean or Romantic tragedy. Such Shakespearean heroes as Macbeth, Richard III, and Mark Antony cannot be regarded as embodiments of any transcendent good. They behave as they do, says Hegel, now speaking outside of his scheme of tragedy, simply because they are the kind of men they are. In a statement pointing up the essence of uninhibited romantic lust and willfulness Hegel said: "it is the inner experience of their heart and individual emotion, or the particular qualities of their personality, which insist on satisfaction."

The traditional categories of tragedy are nearly destroyed in the deepened subjectivities of Romanticism of the 19th-century German philosophers, Arthur Schopenhauer and his disciple Friedrich Nietzsche. In Schopenhauer's *Die Welt als Wille und Vorstellung* (1819; *The World as Will and Idea*), much more than the social or ethical order is upturned. In place of God, the good, reason, soul, or heart, Schopenhauer installs the will, as reality's true inner nature, the metaphysical to everything physical in the world. In Schopenhauer, there is no question of a Hegelian struggle to achieve a more comprehensive good. There is rather the strife of will with itself, manifested by fate in the form of chance and error and by the tragic personages themselves. Both fate and men represent one and the same will, which lives and appears in them all. Its individual manifestations, however, in the form of such phenomena as chances, errors, or men, fight against and destroy each other.

Schopenhauer accordingly rejects the idea of poetic justice: "the demand for so-called poetical justice rests on en-

Tragedy as
Hegelian
dialectic

Schopenhauer's
three types

tire misconception of the nature of tragedy, and, indeed, of the nature of the world itself . . . The true sense of tragedy is the deeper insight, that it is not his own individual sins that the hero atones for, but original sin, *i.e.*, the crime of existence itself. . . . Schopenhauer distinguishes three types of tragic representation: (1) "by means of a character of extraordinary wickedness . . . who becomes the author of the misfortune"; (2) "blind fate—*i.e.*, chance and error" (such as the title characters in Shakespeare's *Romeo and Juliet* and "most of the tragedies of the ancients"); and (3) when "characters of ordinary morality . . . are so situated with regard to each other that their position compels them, knowingly and with their eyes open, to do each other the greatest injury, without any one of them being entirely in the wrong" (such as, "to a certain extent," *Hamlet*).

This last kind of tragedy seems to Schopenhauer far to surpass the other two. His reason, almost too grim to record, is that it provides the widest possible play to the destructive manifestations of the will. It brings tragedy, so to speak, closest to home.

Schopenhauer finds tragedy to be the summit of poetical art, because of the greatness of its effect and the difficulty of its achievement. According to Schopenhauer, the egoism of the protagonist is purified by suffering almost to the purity of nihilism. His personal motives become dispersed as his insight into them grows: "the complete knowledge of the nature of the world, which has a quieting effect on the will, produces resignation, the surrender not merely of life, but of the very will to live."

Schopenhauer's description has limited application to tragic denouements in general. In the case of his own archetypal hero, the hero's end seems merely the mirror image of his career, an oblivion of resignation or death that follows an oblivion of violence. Instead of a dialogue between higher and lower worlds of morality or feeling (which take place even in Shakespeare's darkest plays), Schopenhauer posits a succession of states as helpless in knowledge as in blindness. His "will" becomes a synonym for all that is possessed and necessity-ridden.

Nietzsche's
division
of tragedy
in two
elements

Nietzsche's *Geburt der Tragödie aus dem Geiste der Musik* (1872) was deeply influenced by Schopenhauer. The two elements of tragedy, says Nietzsche, are the Apollonian (related to the Greek god Apollo, here used as a symbol of measured restraint) and the Dionysian (from Dionysus, the Greek god of ecstasy). His conception of the Apollonian is the equivalent of what Schopenhauer called the individual phenomenon—the particular chance, error, or man, the individuality of which is merely a mask for the essential truth of reality which it conceals. The Dionysian element is a sense of universal reality, which, according to Schopenhauer, is experienced after the loss of individual egoism. The "Dionysian ecstasy," as defined by Nietzsche, is experienced "not as individuals but as the *one* living being, with whose creative joy we are united."

Nietzsche dismisses out of hand one of the most venerable features of the criticism of tragedy, the attempt to reconcile the claims of ethics and art. He says that the events of a tragedy are "supposed" to discharge pity and fear and are "supposed" to elevate and inspire by the triumph of noble principles at the sacrifice of the hero. But art, he says, must demand purity within its own sphere. To explain tragic myth, the first requirement is to seek the pleasure that is peculiar to it in the purely aesthetic sphere, without bringing in pity, fear, or the morally sublime.

The essence of this specifically aesthetic tragic effect is that it both reveals and conceals, causing both pain and joy. The drama's exhibition of the phenomena of suffering individuals (Apollonian elements) forces upon the audience "the struggle, the pain, the destruction of phenomena," which in turn communicates "the exuberant fertility of the universal." The spectators then "become, as it were, one with the infinite primordial joy in existence, and . . . we anticipate, in Dionysian ecstasy, the indestructibility and eternity of this joy." Thus, he says, there is a desire "to see tragedy and at the same time to get beyond all seeing . . . to hear and at the same time long to get beyond all hearing."

The inspired force of Nietzsche's vision is mingled with a sense of nihilism:

"only after the spirit of science has been pursued to its limits, . . . may we hope for a rebirth of tragedy . . . I understand by the spirit of science—the faith that first came to light in the person of Socrates—the faith in the explicability of nature and in knowledge as a panacea."

Nietzsche would replace the spirit of science with a conception of existence and the world as an aesthetic phenomenon and justified only as such. Tragedy would enjoy a prominent propagandistic place. It is "precisely the tragic myth that has to convince us that even the ugly and disharmonic are part of an artistic game that the will in the eternal amplitude of its pleasure plays with itself." And, consummately: "we have art in order that we may not perish through truth."

Tragedy in music. Musical dissonance was Nietzsche's model for the double effect of tragedy. The first edition of his book was titled *The Birth of Tragedy out of the Spirit of Music*, another influence from Schopenhauer, for whom music differed from all the other arts in that it is not a copy of a phenomenon but the direct copy of the will itself. He even called the world "embodied music, . . . embodied will." Nietzsche's theorizing on the relation of the tragic theme to art forms other than the drama was in fact confirmed in such operas as Mussorgsky's version of Pushkin's tragedy *Boris Godunov*, Verdi's of *Macbeth* and *Othello*, and Gounod's *Faust*. In contrast to these resettings of received forms, Wagner, Verdi, and Bizet achieved a new kind of tragic power for Romanticism in the theme of the operatic love-death in, respectively, *Tristan and Isolde*, *Aida*, and *Carmen*. Thus, the previous progression of the genre from tragedy to tragicomedy to romantic tragedy continued to a literary-musical embodiment of what Nietzsche called "tragic dithyrambs."

An earlier prophecy than Nietzsche's regarding tragedy and opera was made by the German poet Friedrich von Schiller in a letter of 1797 to Goethe:

I have always trusted that out of opera, as out of the choruses of the ancient festival of Bacchus, tragedy would liberate itself and develop in a nobler form. In opera, servile imitation of nature is dispensed with . . . here is . . . the avenue by which the ideal can steal its way back into the theatre.

20th-century critical theory. In the 20th century, discussion of tragedy was sporadic until the aftermath of World War II. Then it enjoyed new vigour, perhaps to compensate for, or help explain, the dearth of genuine tragic literature, either in the novel or in the theatre. In the 1950s and 1960s countless full-length studies, articles, and monographs variously sought the essence, the vision, the view of life, or the spirit of tragedy out of a concern for the vital culture loss were the death of tragedy to become a reality. They also attempted to mediate the meaning of tragedy to a public that was denied its reality, save in revivals or an occasional approximation. Since the Romantic critics first ventured beyond the Aristotelean categories to consider tragedy, or the tragic, as a sense of life, there was an increasing tendency to regard tragedy not merely as drama but as a philosophical form. It is noteworthy that the Spanish philosopher Miguel de Unamuno's influential book, *The Tragic Sense of Life* (1921), barely mentions the formal drama.

From the time of Aristotle, tragedy has achieved importance primarily as a medium of self-discovery—the discovery of man's place in the universe and in society. That is the main concern of Aristotle in his statements about reversal, recognition, and catharsis, though it remained for the Romantic critics to point it out. The loss of this concern in the facile plays of the 19th and 20th centuries resulted in the reduction of tragic mystery to confused sentimentalism. Critics of the 20th century, being less certain even than Schopenhauer or Nietzsche of what man's place in the scheme of things may be, experimented with a variety of critical approaches, just as contemporary dramatists experimented with various "theatres." Although these critics lacked the philosophical certainties of earlier theorists, they had a richer variety of cultures and genres to instruct them. The hope of both critics and dramatists was that this multiplicity would produce not mere impressionism or haphazard eclecticism but new form and new meaning.

(L.C.)

OTHER GENRES

Satire

"Satire" is a protean term. Together with its derivatives, it is one of the most heavily worked literary designations and one of the most imprecise. The great English lexicographer Samuel Johnson defined satire as "a poem in which wickedness or folly is censured," and more elaborate definitions are rarely more satisfactory. No strict definition can encompass the complexity of a word that signifies, on one hand, a kind of literature—as when one speaks of the satires of the Roman poet Horace or calls the American novelist Nathanael West's *A Cool Million* a satire—and, on the other, a mocking spirit or tone that manifests itself in many literary genres but can also enter into almost any kind of human communication. Wherever wit is employed to expose something foolish or vicious to criticism, there satire exists, whether it be in song or sermon, in painting or political debate, on television or in the movies. In this sense satire is everywhere. Although this section deals primarily with satire as a literary phenomenon, it records its manifestations in a number of other areas of human activity.

THE NATURE OF SATIRE

Historical definitions. The terminological difficulty is pointed up by a phrase of the Roman rhetorician Quintilian: "satire is wholly our own" ("satira tota nostra est"). Quintilian seems to be claiming satire as a Roman phenomenon, although he had read the Greek dramatist Aristophanes and was familiar with a number of Greek forms that one would call satiric. But the Greeks had no specific word for satire; and by *satira* (which meant originally something like "medley" or "miscellany" and from which comes the English "satire") Quintilian intended to specify that kind of poem "invented" by Lucilius, written in hexameters on certain appropriate themes, and characterized by a Lucilian-Horatian tone. *Satura* referred, in short, to a poetic form, established and fixed by Roman practice. (Quintilian mentions also an even older kind of satire written in prose by Marcus Terentius Varro and, one might add, by Menippus and his followers Lucian and Petronius.) After Quintilian's day *satira* began to be used metaphorically to designate works that were "satirical" in tone but not in form. As soon as a noun enters the domain of metaphor, as one modern scholar has pointed out, it clamors for extension; and *satira* (which had had no verbal, adverbial, or adjectival forms) was immediately broadened by appropriation from the Greek word for "satyr" (*satyros*) and its derivatives. The odd result is that the English "satire" comes from the Latin *satira*; but "satirize," "satiric," etc., are of Greek origin. By about the 4th century AD the writer of satires came to be known as *satyricus*; St. Jerome, for example, was called by one of his enemies "a satirist in prose" ("satyricus scriptor in prosa"). Subsequent orthographic modifications obscured the Latin origin of the word satire: *satira* becomes *satyra*, and in England, by the 16th century, it was written "satyre."

Elizabethan writers, anxious to follow Classical models but misled by a false etymology, believed that "satyre" derived from the Greek *satyr* play: satyrs being notoriously rude, unmannerly creatures, it seemed to follow that "satyre" should be harsh, coarse, rough. The English author Joseph Hall wrote:

The Satyre should be like the Poreupine,
That shoots sharpe quills out in each angry line,
And wounds the blushing cheek, and fiery eye,
Of him that heares, and readeth guiltily.

(*Virgidemiarum*, V, 3, 1-4.)

The false etymology that derives satire from satyrs was finally exposed in the 17th century by the Classical scholar Isaac Casaubon; but the old tradition has aesthetic if not etymological appropriateness and has remained strong.

In the prologue to his book, Hall makes a claim that has caused confusion like that following from Quintilian's remark on Roman satire. Hall boasts:

I first adventure: follow me who list,
And be the second English Satyrist.

But Hall knew the satirical poems of Geoffrey Chaucer and John Skelton, among other predecessors, and probably meant that he was the first to imitate systematically the formal satirists of Rome.

Influence of Horace and Juvenal. By their practice, the great Roman poets Horace and Juvenal set indelibly the lineaments of the genre known as the formal verse satire and, in so doing, exerted pervasive, if often indirect, influence on all subsequent literary satire. They gave laws to the form they established, but it must be said that the laws were very loose indeed. Consider, for example, style. In three of his Satires (I, iv; I, x; II, i) Horace discusses the tone appropriate to the satirist who out of a moral concern attacks the vice and folly he sees around him. As opposed to the harshness of Lucilius, Horace opts for mild mockery and playful wit as the means most effective for his ends. Although I portray examples of folly, he says, I am not a prosecutor and I do not like to give pain; if I laugh at the nonsense I see about me, I am not motivated by malice. The satirist's verse, he implies, should reflect this attitude: it should be easy and unpretentious, sharp when necessary, but flexible enough to vary from grave to gay. In short, the character of the satirist as projected by Horace is that of an urbane man of the world, concerned about folly, which he sees everywhere, but moved to laughter rather than rage.

Juvenal, over a century later, conceives the satirist's role differently. His most characteristic posture is that of the upright man who looks with horror on the corruptions of his time, his heart consumed with anger and frustration. Why does he write satire? Because tragedy and epic are irrelevant to his age. Viciousness and corruption so dominate Roman life that for an honest man it is difficult *not* to write satire. He looks about him, and his heart burns dry with rage; never has vice been more triumphant. How can he be silent (*Satires*, I)? Juvenal's declamatory manner, the amplification and luxuriousness of his invective, are wholly out of keeping with the stylistic prescriptions set by Horace. At the end of the scabrous sixth satire, a long, perfervid invective against women, Juvenal flaunts his innovation: in this poem, he says, satire has gone beyond the limits established by his predecessors; it has taken to itself the lofty tone of tragedy.

The results of Juvenal's innovation have been highly confusing for literary history. What *is* satire if the two poets universally acknowledged to be supreme masters of the form differ so completely in their work as to be almost incommensurable? The formulation of the English poet John Dryden has been widely accepted. Roman satire has two kinds, he says: comical satire and tragical satire, each with its own kind of legitimacy. These denominations have come to mark the boundaries of the satiric spectrum, whether reference is to poetry or prose or to some form of satiric expression in another medium. At the Horatian end of the spectrum, satire merges imperceptibly into comedy, which has an abiding interest in the follies of men but has not satire's reforming intent. The distinction between the two modes, rarely clear, is marked by the intensity with which folly is pursued: fops and fools and pedants appear in both, but only satire tries to mend men through them. And, although the great engine of both comedy and satire is irony, in satire, as the 20th-century critic Northrop Frye has said, irony is militant.

Boileau, Dryden, and Alexander Pope, writing in the 17th and 18th centuries—the modern age of satire—catch beautifully, when they like, the deft Horatian tone; however, satire's wit can also be sombre, deeply probing, and prophetic, as it explores the ranges of the Juvenalian end of the satiric spectrum, where satire merges with tragedy, melodrama, and nightmare. Pope's *Dunciad* ends with these lines:

Lo! thy dread Empire, CHAOS! is restor'd;
Light dies before thy uncreating word:

Horace's
concept of
the satirist

Comical
satire and
tragical
satire

Satire and
satyr

Thy hand, great Anarch! lets the curtain fall;
And Universal Darkness buries All.

It is the same darkness that falls on Book IV of Jonathan Swift's *Gulliver's Travels*, on some of Mark Twain's satire—*The Mysterious Stranger, To The Person Sitting in Darkness*—and on George Orwell's 1984.

Structure of verse satire. Roman satire is hardly more determinate in its structure than in its style; the poems are so haphazardly organized, so randomly individual, that there seems little justification for speaking of them as a literary kind at all. Beneath the surface complexity of the poems, however, there exists, as one modern scholar has pointed out, a structural principle common to the satires of the Roman poets and their French and English followers. These poems have a bipartite structure: a thesis part, in which some vice or folly is examined critically from many different angles and points of view; and an antithesis part, in which an opposing virtue is recommended. The two parts are disproportionate in length and in importance, for satirists have always been more disposed to castigate wickedness than exhort to virtue.

Most verse satires are enclosed by a "frame." Just as a novel by the early-20th-century writer Joseph Conrad may be framed by a situation in which his narrator sits on a veranda in the tropics, telling his tale, stimulated into elaboration by the queries of his listeners, so the satire will be framed by a conflict of sorts between the satirist (or, more reasonably, his persona, a fictive counterpart, the "I" of the poem) and an adversary. Usually the adversary has a minor role, serving only to prod the speaker into extended comment on the issue (vice or folly) at hand; he may be sketchily defined, or he may be as effectively projected as Horace's Trebatius (*Satires*, II, i) or his awful bore (I, vi) or his slave Davus, who turns the tables on his master (II, vii). Similarly, the background against which the two talk may be barely suggested, or it may form an integral part of the poem, as in Horace's "Journey to Brundisium" (I, v) or in Juvenal's description of the valley of Egeria, where Umbricius unforgettably pictures the turbulence and decadence of Rome (*Satires*, III). In any event, the frame is usually there, providing a semidramatic situation in which vice and folly may reasonably be dissected.

The satirist has at his disposal an immense variety of literary and rhetorical devices: he may use beast fables, dramatic incidents, fictional experiences, imaginary voyages, character sketches, anecdotes, proverbs, homilies; he may employ invective, sarcasm, burlesque, irony, mockery, raillery, parody, exaggeration, understatement—wit in any of its forms—anything to make the object of attack abhorrent or ridiculous. Amid all this confusing variety, however, there is pressure toward order—internally, from the arraignment of vice and appeal to virtue, and externally, from the often shadowy dramatic situation that frames the poem.

The satiric spirit. Thus, although the formal verse satire of Rome is quantitatively a small body of work, it contains most of the elements later literary satirists employ. When satire is spoken of today, however, there is usually no sense of formal specification whatever; one has in mind a work imbued with the satiric spirit—a spirit that appears (whether as mockery, raillery, ridicule, or formalized invective) in the literature or folklore of all peoples, early and late, preliterate and civilized. According to Aristotle (*Poetics*, IV, 1448b–1449a), Greek Old Comedy developed out of ritualistic ridicule and invective—out of satiric utterances, that is, improvised and hurled at individuals by the leaders of the phallic songs. The function of these "iambic" utterances, it has been shown, was magical; they were thought to drive away evil influences so that the positive fertility magic of the phallus might be operative. This early connection of primitive "satire" with magic has a remarkably widespread history.

In the 7th century BC, the poet Archilochus, said to be the "first" Greek literary satirist, composed verses of such potency against his prospective father-in-law, Lycambes, that Lycambes and his daughter hanged themselves. In the next century the sculptors Bupalus and Athenis "knit their necks in halters," it is said, as a result of the "biting rimes and biting libels" issued by the satirical poet

Hipponax. Similar tales exist in other cultures. The chief function of the ancient Arabic poet was to compose satire (*hijā'*) against the tribal enemy. The satires were thought always to be fatal, and the poet led his people into battle, hurling his verses as he would hurl a spear. Old Irish literature is laced with accounts of the extraordinary power of the poets, whose satires brought disgrace and death to their victims:

... saith [King] Lugh to his poet, "what power can you wield in battle?"

"Not hard to say," quoth Carpre. . . . "I will satirize them, so that through the spell of my art they will not resist warriors."

("The Second Battle of Moytura," trans. by W. Stokes. *Revue Celtique*, XII [1891], 52–130.)

According to saga, when the Irish poet uttered a satire against his victim, three blisters would appear on the victim's cheek, and he would die of shame. One story will serve as illustration: after Deirdriu of the Sorrows came to her unhappy end, King Conchobar fell in love again—this time with the lovely Luaine. They were to be married; but, when the great poet Aithirne the Importunate and his two sons (also poets) saw Luaine, they were overcome with desire for her. They went to Luaine and asked her to sleep with them. She refused. The poets threatened to satirize her. And the story says:

The damsel refused to lie with them. So then they made three satires on her, which left three blotches on her cheeks, to wit, Shame and Blemish and Disgrace. . . . Thereafter the damsel died of shame. . . .

("The Wooing of Luaine. . . ." trans. by W. Stokes. *Revue Celtique*, XXIV [1903], 273–85.)

An eminent 20th-century authority on these matters adduces linguistic, thematic, and other evidence to show a functional relation between primitive "satire," such as that of Carpre and Aithirne, and the "real" satire of more sophisticated times. Today, among various preliterate peoples the power of personal satire and ridicule is appalling; among the Ashanti of West Africa, for example, ridicule is (or was recently) feared more than almost any other humanly inflicted punishment, and suicide is frequently resorted to as an escape from its terrors. Primitive satire such as that described above can hardly be spoken of in literary terms; its affiliations are rather with the magical incantation and the curse.

SATIRICAL MEDIA

Literature. When the satiric utterance breaks loose from its background in ritual and magic, as in ancient Greece (when it is free, that is, to develop in response to literary stimuli rather than the "practical" impulses of magic), it is found embodied in an indefinite number of literary forms that profess to convey moral instruction by means of laughter, ridicule, mockery; the satiric spirit proliferates everywhere, adapting itself to whatever mode (verse or prose) seems congenial. Its targets range from one of Pope's dunces to the entire race of man, as in *Satyr Against Mankind* (1675), by John Wilmot, the earl of Rochester, from Erasmus' attack on corruptions in the church to Swift's excoriation of all civilized institutions in *Gulliver's Travels*. Its forms are as varied as its victims: from an anonymous medieval invective against social injustice to the superb wit of Chaucer and the laughter of Rabelais; from the burlesque of Luigi Pulci to the scurrilities of Pietro Aretino and the "black humour" of Lenny Bruce; from the failings of John Marston and the mordancies of Francisco Gómez de Quevedo y Villegas to the bite of Jean de La Fontaine and the great dramatic structures of Ben Jonson and Molière; from an epigram of Martial to the fictions of Nikolay Gogol and of Günter Grass and the satirical utopias of Yevgeny Zamyatin, Aldous Huxley, and Orwell.

It is easy to see how the satiric spirit would combine readily with those forms of prose fiction that deal with the ugly realities of the world but that satire should find congenial a genre such as the fictional utopia seems odd. From the publication of Thomas More's eponymous *Utopia* (1516), however, satire has been an important ingredient of utopian fiction. More drew heavily on the satire of Horace, Juvenal, and Lucian in composing his

Literary
and rhetor-
ical devices

Power of
ancient
satire

Use of
satire in
fictional
utopias

great work. For example, like a poem by Horace, *Utopia* is framed by a dialogue between "Thomas More" (the historical man a character in his own fiction) and a seafaring philosopher named Raphael Hythloday. The two talk throughout a long and memorable day in a garden in Antwerp. "More's" function is to draw Hythloday out and to oppose him on certain issues, notably his defense of the communism he found in the land of Utopia. "More" is the adversary. Hythloday's role is to expound on the institutions of Utopia but also to expose the corruption of contemporary society. Thus he functions as a satirist. Here Hythloday explains why Englishmen, forced off their land to make way for sheep, become thieves:

Forsooth . . . your sheep that were wont to be so meek and tame and so small eaters, now as I hear say, be become so great devourers and so wild, that they eat up and swallow down the very men themselves. They consume, destroy, and devour whole fields, houses, and cities. For look in what parts of the realm doth grow the finest and therefore dearest wool, there noblemen and gentlemen, yea and certain abbots, holy men no doubt, not contenting themselves with the yearly revenues and profits that were wont to grow to their forefathers and predecessors of their lands, nor being content that they live in rest and pleasure nothing profiting, yea, much annoying the weal-public, leave no ground for tillage. They enclose all into pastures; they throw down houses; they pluck down towns and leave nothing standing but only the church to be made a sheep-house.

(More's *Utopia*, Everyman edition, 1951.)

Here are characteristic devices of the satirist, dazzlingly exploited: the beast fable compressed into the grotesque metaphor of the voracious sheep; the reality-destroying language that metamorphoses gentlemen and abbots into earthquakes and a church into a sheep barn; the irony coldly encompassing the passion of the scene. Few satirists of any time could improve on this.

Just as satire is a necessary element of the work that gave the literary form utopia its name, so the utopias of Lilliput, Brobdingnag, and Houyhnhnmland are essential to the satire of More's great follower Jonathan Swift. He sent Gulliver to different lands from those Hythloday discovered, but Gulliver found the same follies and the same vices, and he employed a good many of the same rhetorical techniques his predecessor had used to expose them. *Gulliver's Travels*, as one scholar points out, is a salute across the centuries to Thomas More. With this kind of precedent, it is not surprising that in the 20th century, when utopia turns against itself, as in Aldous Huxley's *Brave New World* (1932), the result is satire unrelieved.

Drama. The drama has provided a favourable environment for satire ever since it was cultivated by Aristophanes, working under the extraordinarily open political conditions of 5th-century Athens. In a whole series of plays—*The Clouds*, *The Frogs*, *Lysistrata*, and many others—Aristophanes lampoons the demagogue Cleon by name, violently attacks Athenian war policy, derides the audience of his plays for their gullible complacency, pokes fun at Socrates as representative of the new philosophical teaching, stages a brilliantly parodic poetic competition between the dramatists Aeschylus and Euripides in Hades, and in general lashes out at contemporary evils with an uninhibited and unrivalled inventiveness. But the theatre has rarely enjoyed the political freedom Aristophanes had—one reason, perhaps, that satire more often appears in drama episodically or in small doses than in the full-blown Aristophanic manner. In Elizabethan England, Ben Jonson wrote plays that he called "comicall satyres"—*Every Man Out of His Humour*, *Poetaster*—and there are substantial elements of satire in Shakespeare's plays—some in the comedies, but more impressively a dark and bitter satire in *Timon of Athens*, *Troilus and Cressida*, *Hamlet*, and *King Lear*. The 17th-century comedy of Molière sometimes deepens into satire, as with the exposure of religious hypocrisy in *Tartuffe* or the railing against social hypocrisy by Alceste in *The Misanthrope*. George Bernard Shaw considered himself a satirist. He once compared his country's morals to decayed teeth and himself to a dentist, obliged by his profession to give pain in the interests of better health. Yet, as inventive and witty as Shaw is,

compared to the 20th-century German playwright Bertolt Brecht, whose anatomizing of social injustice cuts deep, Shaw is a gentle practitioner indeed.

Motion pictures and television. The movies have sometimes done better by satire than the theatre, and it is in the movies that an ancient doctrine having to do with principles of decorum in the use of satire and ridicule has been exploded. The English novelist Henry Fielding was reflecting centuries of tradition when, in the preface to *Joseph Andrews* (1742), he spoke of the inappropriateness of ridicule applied to black villainy or dire calamity. "What could exceed the absurdity of an Author, who should write *the Comedy of Nero, with the merry Incident of ripping up his Mother's Belly?*" Given this point of view, Hitler seems an unlikely target for satire; yet in *The Great Dictator* (1940) Charlie Chaplin managed a successful, if risky, burlesque. Chaplin has written, however, that, determined as he was to ridicule the Nazi notions of a superrace, if he had known of the horrors of the concentration camps, he could not have made the film. Stanley Kubrick's *Dr. Strangelove; or, How I Learned to Stop Worrying and Love the Bomb* (1964) denies all limitation; through some alchemy Kubrick created an immensely funny, savagely satirical film about the annihilation of the world. A combination of farce and nightmare, *Dr. Strangelove* satirizes military men, scientists, statesmen—the whole ethos of the technological age—in the most mordant terms; it shows the doomsday blast, yet leaves audiences laughing. "You can't fight in here," says the president of the United States as doom nears, "this is the War Room." The film's tone is less didactic than in most powerful satire—the mushroom cloud carries its own moral—yet satire's full force is there.

Television has not proved a notably receptive medium. *That Was the Week That Was*, a weekly satirical review started in England in 1962, had remarkable success for a time but succumbed to a variety of pressures, some of them political; when a version of the program was attempted in the United States, it was emasculated by restrictions imposed by sponsors fearful of offending customers and by program lawyers wary of libel suits. Jonathan Swift said that he wrote to vex the world rather than divert it; it is not an attitude calculated to sell soap.

Festivals. Yet satire does much more than vex, and even in Swift's work there is a kind of gaiety that is found in many nonliterary manifestations of the satirical spirit. Satire always accompanies certain festivals, for example, particularly saturnalian festivals. Many different cultures set aside a holiday period in which customary social restraints are abandoned, distinctions of rank and status are turned upside down, and institutions normally sacrosanct are subjected to ridicule, mockery, burlesque. The Romans had their Saturnalia, the Middle Ages its Feast of Fools; and in the 20th century many countries still have annual carnivals (Fasching in Austria, the Schnitzelbank in Basel, Switzerland, for example) at which, amid other kinds of abandon, an extraordinary freedom of satirical utterance is permitted. Even in Africa among the Ashanti, for whom ridicule has such terrors, there is a festival during which the sacred chief himself is satirized. "Wait until Friday," said the chief to the enquiring anthropologist, "when the people really begin to abuse me, and if you will come and do so too it will please me." Festivals such as these provide sanctioned release from social inhibition and repression, and, in these circumstances, satire directed at men in power or at taboo institutions acts as a safety valve for pent-up frustrations.

Satire may often function this way. A story is told that the 16th-century pope Adrian VI was highly offended at satirical verses written against him and affixed to Pasquino's statue (a famous repository for lampoons in Rome), but he became a willing target once he realized that his enemies vented otherwise dangerous hostility in this relatively harmless manner. Similar mechanisms operate today when, at a nightclub or theatre, audiences listen to satirical attacks on political figures or on issues such as racial discrimination, identify with the satirist, laugh at his wit, and thereby discharge their own aggressive feelings. Satire of this order

Aristophanes' use of satire

The gaiety of satire

Caricature and cartoon

is a far cry, of course, from that written by a Swift or Voltaire, whose work can be said to have a revolutionary effect.

Visual arts. The critique of satire may be conveyed even more potently in the visual arts than by way of the spoken or written word. In caricature and in what came to be known as the cartoon, artists since the Renaissance have left a wealth of startlingly vivid commentary on the men and events of their time. The names alone evoke their achievement: in England, William Hogarth, Thomas Rowlandson, Sir John Tenniel, and Sir Max Beerbohm; in France, Charles Philipon (whose slow-motion metamorphosis of King Louis-Philippe into a *poire*—that is, “fathead,” or “fool”—is classic) and Honoré Daumier; in Spain, Francisco Goya, and out of Spain, Pablo Picasso; and among recent political cartoonists, Sir David Low, Vicky (Victor Welsz), Herblock (Herbert Block), and Conrad.

The favourite medium of such individuals is the black-and-white print in which the satirical attack is pointed up by a brief verbal caption. The social impact of their art is incalculable. Dictators recognize this all too well, and in times of social tension political cartoonists are among the first victims of the censor.

THE SATIRIST, THE LAW, AND SOCIETY

Indeed, the relations of satirists to the law have always been delicate and complex. Both Horace and Juvenal took extraordinary pains to avoid entanglements with authority—Juvenal ends his first satire with the self-protective announcement that he will write only of the dead. In England in 1599 the Archbishop of Canterbury and the Bishop of London issued an order prohibiting the printing of any satires whatever and requiring that the published satires of Hall, John Marston, Thomas Nashe, and others be burned.

Today the satirist attacks individuals only at the risk of severe financial loss to himself and his publisher. In totalitarian countries he even risks imprisonment or death. Under extreme conditions satire against the reigning order is out of the question. Such was the case in the Soviet Union and most other Communist countries. For example, the poet Osip Mandelstam was sent to a concentration camp and his death for composing a satirical poem on Stalin.

One creative response the satirist makes to social and legal pressures is to try by rhetorical means to approach his target indirectly; that is, a prohibition of direct attack fosters the manoeuvres of indirection that will make the attack palatable: *e.g.*, irony, burlesque, and parody. It is a nice complication that the devices that render satire acceptable to society at the same time sharpen its point. “Abuse is not so dangerous,” said Dr. Johnson, “when there is no vehicle of wit or delicacy, no subtle conveyance.” The conveyances are born out of prohibition.

Anthony Cooper, 3rd earl of Shaftesbury, writing in the 18th century, recognized the “creative” significance of legal and other repressions on the writing of satire. “The greater the weight [of constraint] is, the bitterer will be the satire. The higher the slavery, the more exquisite the buffoonery.” Shaftesbury’s insight requires the qualification made above. Under a massively efficient tyranny, satire of the forms, institutions, or personalities of that tyranny is impossible. But, under the more relaxed authoritarianism of an easier going day, remarkable things could be done. Max Radin, a Polish-born American author, noted how satirical journals in Germany before World War I, even in the face of a severe law, vied with each other to see how close they could come to caricatures of the Kaiser without actually producing them. “Satire which the censor understands,” said the Austrian satirist Karl Kraus, “deserves to be banned.”

The 20th-century American critic Kenneth Burke summed up this paradoxical aspect of satire’s relation with the law by suggesting that the most inventive satire is produced when the satirist knowingly takes serious risks and is not sure whether he will be acclaimed or punished. The whole career of Voltaire is an excellent case in point. Bigots and tyrants may have turned pale at his name, as a famous hyperbole has it; however, Voltaire’s satire was

sharpened and his life rendered painfully complicated as he sought to avoid the penalties of the law and the wrath of those he had angered. Men such as Voltaire and Kraus and the Russian Ye.I. Zamyatin attack evil in high places, pitting their wit and moral authority against cruder forms of power. In this engagement there is frequently something of the heroic.

Readers have an excellent opportunity to examine the satirist’s claim to social approval by reason of the literary convention that decrees that he must justify his problematic art. Nearly all satirists write apologies, and nearly all the apologies project an image of the satirist as a plain, honest man, wishing harm to no worthy person but appalled at the evil he sees around him and forced by his conscience to write satire. Pope’s claim is the most extravagant:

Yes, I am proud; I must be proud to see
Men not afraid of God, afraid of me;
Safe from the Bar, the Pulpit, and the Throne,
Yet touch’d and sham’d by *Ridicule* alone.
O sacred Weapon! left for Truth’s defence,
Sole Dread of Folly, Vice, and Insolence!

(*Epilogue to the Satires*, II, 208–13.)

After the great age of satire, which Pope brought to a close, such pretensions would have been wholly anachronistic. Ridicule depends on shared assumptions against which the deviant stands in naked relief. The satirist must have an audience that shares with him commitment to certain intellectual and moral standards that validate his attacks on aberration. The greatest satire has been written in periods when ethical and rational norms were sufficiently powerful to attract widespread assent yet not so powerful as to compel absolute conformity—those periods when the satirist could be of his society and yet apart, could exercise a double vision.

Neoclassic writers had available to them as an implicit metaphor the towering standard of the classical past; for the 19th and 20th centuries no such metaphors have been available. It is odd, however, that, whereas the 19th century in general disliked and distrusted satire (there are of course obvious exceptions), our own age, bereft of unifying symbols, scorning traditional rituals, searching for beliefs, still finds satire a congenial mode in almost any medium. Although much of today’s satire is self-serving and trivial, there are notable achievements. Joseph Heller’s novel *Catch-22* (1961) once again makes use of farce as the agent of the most probing criticism: Who is sane, the book asks, in a world whose major energies are devoted to blowing itself up? Beneath a surface of hilariously grotesque fantasy, in which characters from Marx brothers’ comedy carry out lethal assignments, there is exposed a dehumanized world of hypocrisy, greed, and cant. Heller is a satirist in the great tradition. If he can no longer, like Pope, tell men with confidence what they should be for, he is splendid at showing them what they must be against. The reader laughs at the mad logic of *Catch-22*, and, as he laughs, he learns. This is precisely the way satire has worked from the beginning. (R.C.E.)

Nonfictional prose

Defining nonfictional prose literature is an immensely challenging task. Nonfictional prose literature differs from bald statements of fact, such as those recorded in an old chronicle or inserted in a business letter or in an impersonal message of mere information. As used in a broad sense, the term nonfictional prose literature here designates writing intended to instruct (but not highly scientific and erudite writings in which no aesthetic concern is evinced), to impart wisdom or faith, and especially to please. Separate sections cover biographical literature and literary criticism.

NATURE

Nonfictional prose genres cover an almost infinite variety of themes, and they assume many shapes. In quantitative terms, if such could ever be valid in such nonmeasurable matters, they probably include more than half of all that has been written in countries having a literature of their

Relationship between the satirist and his audience

Repressions of satire

own. Nonfictional prose genres have flourished in nearly all countries with advanced literatures. The genres include political and polemical writings, biographical and autobiographical literature, religious writings, and philosophical, and moral or religious writings.

After the Renaissance, from the 16th century onward in Europe, a personal manner of writing grew in importance. The author strove for more or less disguised self-revelation and introspective analysis, often in the form of letters, private diaries, and confessions. Also of increasing importance were aphorisms after the style of the ancient Roman philosophers Seneca and Epictetus, imaginary dialogues, and historical narratives, and later, journalistic articles and extremely diverse essays. From the 19th century, writers in Romance and Slavic languages especially, and to a far lesser extent British and American writers, developed the attitude that a literature is most truly modern when it acquires a marked degree of self-awareness and obstinately reflects on its purpose and technique. Such writers were not content with imaginative creation alone; they also explained their work and defined their method in prefaces, reflections, essays, self-portraits, and critical articles. The 19th-century French poet Charles Baudelaire asserted that no great poet could ever quite resist the temptation to become also a critic: a critic of others and of himself. Indeed, most modern writers, in lands other than the United States, whether they be poets, novelists, or dramatists, have composed more nonfictional prose than poetry, fiction, or drama. In the instances of such monumental figures of 20th-century literature as the poets Ezra Pound, T.S. Eliot, and perhaps William Butler Yeats, or the novelists Thomas Mann and André Gide, that part of their output may well be considered by posterity to be equal in importance to their more imaginative writing.

It is virtually impossible to attempt a unitary characterization of nonfictional prose. The concern that any definition is a limitation, and perhaps an exclusion of the essential, is nowhere more apposite than to this inordinately vast and variegated literature. Ever since the ancient Greek and Roman philosophers devised literary genres, some critics have found it convenient to arrange literary production into kinds or to refer it to modes.

ELEMENTS

Obviously, a realm as boundless and diverse as nonfictional prose literature cannot be characterized as having any unity of intent, of technique, or of style. It can be defined, very loosely, only by what it is not. Many exceptions, in such a mass of writings, can always be brought up to contradict any rule or generalization. No prescriptive treatment is acceptable for the writing of essays, of aphorisms, of literary journalism, of polemical controversy, of travel literature, of memoirs and intimate diaries. No norms are recognized to determine whether a dialogue, a confession, a piece of religious or of scientific writing, is excellent, mediocre, or outright bad, and each author has to be relished, and appraised, chiefly in his own right. "The only technique," the English critic F.R. Leavis wrote in 1957, "is that which compels words to express an intensively personal way of feeling." Intensity is probably useful as a standard; yet it is a variable, and often elusive, quality, possessed by polemicists and by ardent essayists to a greater extent than by others who are equally great. "Loving, and taking the liberties of a lover" was Virginia Woolf's characterization of the 19th-century critic William Hazlitt's style; it instilled passion into his critical essays. But other equally significant English essayists of the same century, such as Charles Lamb or Walter Pater, or the French critic Hippolyte Taine, under an impassive mask, loved too, but differently. Still other nonfictional writers have been detached, seemingly aloof, or, like the 17th-century French epigrammatist La Rochefoucauld, sarcastic. Their intensity is of another sort.

Reality and imagination. Prose that is nonfictional is generally supposed to cling to reality more closely than that which invents stories, or frames imaginary plots. Calling it "realistic," however, would be a gross distortion. Since nonfictional prose does not stress inventiveness of themes and of characters independent of the author's self,

it appears in the eyes of some moderns to be inferior to works of imagination. In the middle of the 20th century an immensely high evaluation was placed on the imagination, and the adjective "imaginative" became a grossly abused cliché. Many modern novels and plays, however, were woefully deficient in imaginative force, and the word may have been bandied about so much out of a desire for what was least possessed. Many readers are engrossed by travel books, by descriptions of exotic animal life, by essays on the psychology of other nations, by Rilke's notebooks or by Samuel Pepys's diary far more than by poetry or by novels that fail to impose any suspension of disbelief. There is much truth in Oscar Wilde's remark that "the highest criticism is more creative than creation and the primary aim of the critic is to see the object as in itself it really is not." A good deal of imagination has gone not only into criticism but also into the writing of history, of essays, of travel books, and even of the biographies or the confessions that purport to be true to life as it really happened, as it was really experienced.

The imagination at work in nonfictional prose, however, would hardly deserve the august name of "primary imagination" reserved by the 19th-century English poet Samuel Taylor Coleridge to creators who come close to possessing semidivine powers. Rather, imagination is displayed in nonfictional prose in the fanciful invention of decorative details, in digressions practiced as an art and assuming a character of pleasant nonchalance, in establishing a familiar contact with the reader through wit and humour. The variety of themes that may be touched upon in that prose is almost infinite. The treatment of issues may be ponderously didactic and still belong within the literary domain. For centuries, in many nations, in Asiatic languages, in medieval Latin, in the writings of the humanists of the Renaissance, and in those of the Enlightenment, a considerable part of literature has been didactic. The concept of art for art's sake is a late and rather artificial development in the history of culture, and it did not reign supreme even in the few countries in which it was expounded in the 19th century. The ease with which digressions may be inserted in that type of prose affords nonfictional literature a freedom denied to writing falling within other genres. The drawback of such a nondescript literature lies in judging it against any standard of perfection, since perfection implies some conformity with implicit rules and the presence, however vague, of standards such as have been formulated for comedy, tragedy, the ode, the short story and even (in this case, more honoured in the breach than the observance) the novel. The compensating grace is that in much nonfictional literature that repudiates or ignores structure the reader is often delighted with an air of ease and of nonchalance and with that rarest of all virtues in the art of writing: naturalness.

Style. The writing of nonfictional prose should not entail the tension, the monotony, and the self-conscious craft of fiction writing. The search for *le mot juste* ("the precise word") so fanatically pursued by admirers of Flaubert and Maupassant is far less important in nonfictional prose than in the novel and the short story. The English author G.K. Chesterton (1874–1936), who was himself more successful in his rambling volumes of reflections and of religious apologetics than in his novels, defined literature as that rare, almost miraculous use of language "by which a man really says what he means." In essays, letters, reporting, and narratives of travels, the author's aim is often not to overpower his readers by giving them the impression that he knows exactly where he is leading them, as a dramatist or a detective-story writer does. Some rambling casualness, apparently irrelevant anecdotes, and suggestions of the conclusions that the author wishes his readers to infer are often more effective than extreme terseness.

There is also another manner of writing that is more attentive to the periodic cadences and elegance of prose, in the style of the ancient Roman orator Cicero. The 19th-century English essayist William Hazlitt praised the felicities of style and the refinements of the prose of the British statesman Edmund Burke (1729–97) as "that which went the nearest to the verge of poetry and yet never fell over." A number of English writers have been fond of

Imagination in nonfictional prose

The role of eloquence

that harmonious, and rhetorical prose, the taste for which may well have been fostered not only by the familiarity with Cicero but also by the profound influence of the authorized version of the Bible (1611). Martin Luther's translation of the New Testament (1522) and of the Old Testament (1534) likewise molded much of German prose and German sensibility for centuries.

In the 20th century that type of prose lost favour with American and British readers, who ceased to cherish Latin orators and Biblical prose as their models. In German literature, however, in which harmonious balance and eloquence were more likely to be admired, and in other languages more directly derived from Latin, a musical style, akin to a prolonged poem in prose, was cultivated more assiduously, as exemplified in Italian in the writings of Gabriele D'Annunzio, in French in those by André Gide, and in German in *Die Aufzeichnungen des Malte Laurids Brügge* (*The Notebooks of Malte Laurids Brügge*) by the poet Rainer Maria Rilke. Such an elaborate style appears to be more easily tolerated by the readers in non-fictional writing, with its lack of cumulative continuity and, generally speaking, its more restricted size, than in novels such as Pater's *Marius the Epicurean* (1885) and occasionally in Thomas Mann's fiction, in which such a style tends to pall on the reader. Similarly, it is easier for the non-fictional prose writer to weave into his style faint suggestions of irony, archaisms, alliterations, and even interventions of the author that might prove catastrophic to credibility in fiction. Critics have argued that too close attention to style was harmful to the sweep necessary to fiction: they have contended that many of the greatest novelists, such as Dickens, Balzac, Dostoyevsky, and Zola at times "wrote" badly; assuredly, they treated language carelessly more than once. Essayists, historians, orators, and divines often affect a happy-go-lucky ease so as to put them on the same footing with the common reader, but they realize that language and style are vital. They must know what resources they can draw from vivid sensations, brilliant similes, balanced sentences, or sudden, epigrammatic, effects of surprise.

Author presence. The one feature common to most authors of non-fictional prose (a few staid historians and even fewer philosophers excepted) is the marked degree of the author's presence in all they write. That is to be expected in epistolary literature, and, although less inevitably, in the essay, the travel book, journalistic reporting, and polemical or hortatory prose. Although the 17th-century French religious philosopher Pascal hinted that "the ego is hateful," the author's presence is still strongly felt. This presence endows their works with a personal and haunting force that challenges, converts, or repels, but hardly ever leaves the reader indifferent. Saint Paul's epistles owe their impact—perhaps second to none in the history of the Western World—to the self that vehemently expresses itself in them, showing no concern whatever for the niceties of Attic prose. In the treatises, discourses, and philosophical argumentation of the great writers of the Enlightenment, such as Voltaire, Diderot, and Rousseau, they frequently resort to the first person singular, which results in a vivid concreteness in the treatment of ideas. To think the abstract concretely, a precept reminiscent of the 18th-century philosophers, was also the aim of the 20th-century philosophers Jean-Paul Sartre and Maurice Merleau-Ponty when they naturalized Existentialist thought in France. The growth of personal literature in its myriad shapes is one of the striking features of modern literary evolution.

The very fact that the writer of non-fictional prose does not seek an imaginary projection to impart his vision, his anguish, and his delights to readers also underlines the nature of his intention. A school of critics has vigorously attacked "the intentional fallacy," which leads biographers and some literary historians to ask what an artist intended before evaluating the completed work of art. But in a work of apologetics or of homiletics, in a work of history or of sociology, in a critical or even in a desultory and discursive essay, and certainly in aphorisms or maxims or both, the intention of the author remains omnipresent. This intention may be disguised under the mask of a para-

ble, under the interlocutors of a philosophical dialogue, or under the admonitions of a prophet, but the reader is never oblivious of the thinker's intent. The reader has a sophisticated enjoyment of one who shares the creator's intent and travels familiarly along with him. He respects and enjoys in those authors the exercise of an intelligence flexible enough to accept even the irrational as such.

APPROACHES

In terms of approach, that is, the attitude of the writer as it can be inferred from the writing, the distinguishing features of non-fictional prose writings are the degree of presence of the ego and of the use of a subjective, familiar tone. Such devices are also used, of course, by authors of fiction, but to a lesser extent. Similarly, the basic modes of writing—the descriptive, the narrative, the expository, and the argumentative—are found in both non-fictional literature and in fiction, but in different degrees.

The descriptive mode. In non-fictional prose, essayists, moralists, naturalists, and others regularly evoked nature scenes. The most sumptuous masters of prose composed landscapes as elaborately as landscape painters. The French writer and statesman Chateaubriand (1768–1848), for example, who was not outstandingly successful in inventing plots or in creating characters independent from his own self, was a master of description: his writings influenced the French Romantic poets, who set the impassive splendour of outward nature in contrast to the inner anguish of mortals. The 19th-century English art critic John Ruskin had a more precise gift of observation, as revealed in his descriptions of Alpine mountains and of the humblest flowers or mosses, but his ornate and sonorous prose was the climax of a high-flown manner of writing that later read like the majestic relic of another era. American non-fictional writers of the same period such as Ralph Waldo Emerson and Henry Thoreau scrupulously described the lessons of organization, of unity, and of moral beauty to be deciphered from the vicissitudes of nature. Russian essayists vied with novelists in their minute yet rapturous descriptions of the thaw releasing the torrents of spring or the implacable force of the long Northern winters. Writers more inclined to the observation of social life, in satirical sketches of the mechanically polite and artificial habitués of salons, helped the novel of social life come into existence in several Western countries.

Narrative. The narrative element is less conspicuous in writing that does not purport to relate a story than in fictional works, but there is a role for narrative in letters, diaries, autobiographies, and historical writing. Most often, an incident is graphically related by a witness, as in letters or memoirs; an anecdote may serve to illustrate a moral advice in an essay; or an entertaining encounter may be inserted into an essay or a travel sketch. Digression here represents the utmost in art: it provides a relief from the persistent attention required when the author is pursuing his purpose more seriously. Similarly, such writing provides a pleasant contrast to the rigid structure of the majority of novels since the late 19th century. In historical writing, however, simplicity and clarity of narrative are required, though it may be interspersed with speeches, with portraits, or with moral and polemical allusions. In other forms of non-fictional prose, the meandering fancy of the author may well produce an impression of freedom and of truth to life unattainable by the more carefully wrought novel. Many writers have confessed to feeling relieved when they ceased to create novels and shifted to impromptu sketches or desultory essays. The surrealist essayists of the 20th century poured their scorn on detective fiction as the most fiercely logical form of writing. In contrast, the author of essays or other non-fictional prose may blend dreams and facts, ventures into the illogical, and delightful eccentricities.

Expository and argumentative modes. The rules of old-fashioned rhetoric apply better to expository and argumentative prose than to the other modes. These rules were first set down in ancient Greece by teachers who elicited them from the smooth eloquence of Socrates, the impassioned and balanced reasoning of Demosthenes, and others. The ancient Romans went further still in codifying

figures of speech, stylistic devices, and even the gestures of the orator. Such treatises played a significant part in the education of the Renaissance Humanists, of the classical and Augustan prose writers of 17th-century England and France, of the leaders of the French Revolution in the 18th century, and even in 19th-century historians and statesmen such as Guizot in France and Macaulay and Gladstone in Britain. But the sophisticated oratory of such 18th-century British orators as Richard Brinsley Sheridan, Edmund Burke, and Charles Fox or, more recently, that of Winston Churchill, hardly seems attuned to audiences in the age of television.

It has been suggested by students of German history that Adolf Hitler, in his vituperative speeches at Nuremberg in the 1930s, fascinated the Germans because they had been unaccustomed, unlike other Western nations, to eloquence in their leaders. If a large part of a population is illiterate, such as the Cubans under Fidel Castro, unending flows of eloquence may constitute a convenient means of educating the masses. Elsewhere, a more familiar and casual type of address from political leaders tends to be preferred in an era of mass media. The gift of a superior orator has been facetiously defined as that of saying as little as possible in as many words as possible. Like sermons, many types of formal address such as lectures, political speeches, and legal pleadings appear to be doomed as documents of literary value, as Burke's or Lincoln's orations and addresses were when they were learned by heart by the younger generations and helped mold the style and contribute to the moral education of men.

THE ESSAY

In modern literatures, the category of nonfictional prose that probably ranks as the most important both in the quantity and in the quality of its practitioners is the essay.

Modern origins. Before the word itself was coined in the 16th century by Montaigne and Bacon, what came to be called an essay was called a treatise, and its attempt to treat a serious theme with consistency deprived it of the seductive charm relished in the later examples of that form of literature. In this sense, the word "essay" would hardly fit the didactic tone of Aristotle's *Rhetoric* or his *Metaphysics*. There were, however, ancient masters of an early form of the essay, such as Cicero discoursing on the pleasantness of old age or on the art of "divination"; Seneca, on anger or clemency; and Plutarch, more superficially and casually, on the passing of oracles. The relentless desire to analyze one's own contradictions, especially among Christians, who, like Saint Paul, were aware of their duality and of "doing the evil which they would not," also contributed to the emergence of the essay. But Christian writing tended to be highly didactic, as may be seen in the work of Saint Augustine of the 5th century, or of the 12th-century theologian Abelard, or even in the Latin writings on "the solitary life" or on "the scorn of the world" by the 14th-century Italian poet Petrarch. Not until the Renaissance, with its increasing assertion of the self, was the flexible and deliberately nonchalant and versatile form of the essay perfected by Montaigne.

Montaigne, who established the term essay, left his mark on almost every essayist who came after him in continental Europe, and perhaps even more in English-speaking countries. Emerson made him one of his six *Representative Men* along with others of the stature of Plato, Shakespeare, and Goethe. Hazlitt lauded Montaigne's qualities as precisely those that "we consider in great measure English," and another English romantic writer, Leigh Hunt, saw him as "the first man who had the courage to say as an author what he felt as a man." And the 20th-century poet T.S. Eliot declared him to be the most important writer to study for an insight into the literature of France. With Montaigne, the essay achieved for the first time what it can achieve better than any other form of writing, except perhaps the epistolary one: a means of self-discovery. It gave the writer a way of reaching the secret springs of his behaviour, of seizing the man and the author at once in his contradictions, in his profound disunity, and in his mobility. The essay was symbolic of man's new attitude toward himself, revelling in change, and hence in growth,

and forsaking his age-old dream of achieving an underlying steadfastness that might make him invulnerable and similar to the gods. Now he set out to accept himself whole, with his body and his physical and behavioural peculiarities, and thereby repudiate medieval asceticism. He would portray his foibles and unworthiness, hoping to rise above his own mediocrity, or, at the other extreme, he would exalt himself in the hope that he might become the man he depicted. Montaigne in his essays pursued an ethical purpose, but with no pompousness or rhetoric. He offered an ideal that was adopted by his successors for 200 years: perfecting man as a tolerant, undogmatic, urbane social being. But, unlike medieval Christian writers, he would not sacrifice to others the most dearly cherished part of himself. To others he would lend himself, but his personality and his freedom were his own, and his primary duty was to become a wiser human being.

No essayist after Montaigne touched on so many varied aspects of life with such an informal, felicitous, and brilliant style. The later writers who most nearly recall the charm of Montaigne include, in England, Robert Burton, though his whimsicality is more erudite, Sir Thomas Browne, and Laurence Sterne, and in France, with more self-consciousness and pose, André Gide and Jean Cocteau.

Uses of the essay. In the age that followed Montaigne's, at the beginning of the 17th century, social manners, the cultivation of politeness, and the training of an accomplished gentleman became the theme of many essayists. This theme was first exploited by the Italian Baldassare Castiglione in his *Il cortegiano* (1528; *The Courtier*). The influence of the essay and of genres allied to it, such as maxims, portraits, and sketches, proved second to none in molding the behaviour of the cultured classes, first in Italy, then in France, and, through French influence, in most of Europe in the 17th century. Among those who pursued this theme was the 17th-century Spanish Jesuit Baltasar Gracián in his essays on the art of worldly wisdom.

With the advent of a keener political awareness with the age of Enlightenment, in the 18th century, the essay became all-important as the vehicle for a criticism of society and of religion. Because of its flexibility, its brevity, and its potential both for ambiguity and for allusions to current events and conditions, it was an ideal tool for philosophical reformers. *The Federalist Papers* in America and the tracts of the French Revolutionaries, are among the countless examples of attempts during this period to improve the condition of man through the essay.

The advantage of this form of writing was that it was not required to conform to any unity of tone or to similar strictures assigned to other genres since it was for a long time not even considered a genre. After ponderous apologies for traditional faith failed to repulse the onslaught of deism and atheism, traditionalists of the 18th and 19th centuries, such as Burke and Coleridge, abandoned unwieldy dogmatic demonstrations in favour of the short, provocative essay. In the defense of the past, it served as the most potent means of educating the masses. French Catholics, German pietists, and a number of individual English and American authors confided to the essay their dismay at what they saw as modern vulgarity and a breakdown of the coherence of the Western tradition. Essays such as Paul Elmer More's long series of *Shelburne Essays* (published between 1904 and 1935), T.S. Eliot's *After Strange Gods* (1934) and *Notes Towards the Definition of Culture* (1948), and others that attempted to reinterpret and redefine culture, established the genre as the most fitting to express the genteel tradition at odds with the democracy of the new world.

The proliferation of magazines in the United States, and the public's impatience with painstaking demonstrations and polemics, helped establish the essay just as firmly as a receptacle for robust, humorous common sense, unpretentiously expressed, as in the writings of Oliver Wendell Holmes (1809-94). Creative writers resorted to it to admonish their compatriots when they seemed too selfishly unconcerned by the tragedies of the world. Archibald MacLeish, for instance, did so in *A Time to Speak* (1941). Lewis Mumford, Allen Tate, and other literary and social critics became crusaders for moral and spiritual reform;

Montaigne's pre-eminence as an essayist

Journalistic essays

others seized upon the essay for scathingly ironical and destructive criticism of their culture: for example, James Gibbons Huneker (1860–1921), an admirer of iconoclasts and of egoists, as he called them, proposed European examples to Americans he deemed to be too complacent and lethargic; and, more vociferously still, H.L. Mencken (1880–1956), a self-appointed foe of prejudices, substituted his own for those he trounced in his contemporaries.

In other new countries, or in cultures acquiring an awareness of their own ambitious identity, the essay became semipolitical, earnestly nationalistic, and often polemical, playful, or bitter. Such essays sometimes succeeded in shaking the elite out of its passivity. In Uruguay, for example, José Enrique Rodó (1872–1917), in an analogy to the characters in Shakespeare's *Tempest*, compared what should be the authentic South American to the spirit Ariel, in a work thus entitled, in contrast to the bestial Caliban, representing the materialism of North America. In Canada Olivier Asselin (1874–1937) used the essay to advocate the development of a genuine French-Canadian literature. Among the older cultures of Europe, Salvatore Quasimodo (1901–68), the Italian poet and Nobel laureate, appended critical and hortatory essays to some of his volumes of verse, such as *Il falso e vero verde* (1956; "The False and True Green"). Other European heirs to this tradition of the essay include Stefan Zweig and Hugo von Hofmannsthal in Austria and Thomas Mann and Bertolt Brecht in Germany; their sprightly and incisive essays on the arts can be traced to the 19th century German philosopher Arthur Schopenhauer.

One of the functions of literature is to please and to entertain; and the essay, as it grew into the biggest literary domain of all, did not lose the art of providing escape. Essayists have written with grace on children, on women, on love, on sports, as in Robert Louis Stevenson's collection *Virginibus Puerisque* (1881), or Willa Cather's pleasant reflections in *Not Under Forty* (1936). Ernest Renan (1823–92), one of the most accomplished French masters of the essay, found relief from his philosophical and historical studies in his half-ironical considerations on love, and Anatole France (1844–1924), his disciple, and hosts of others have alternated playful essays with others of high seriousness. Sports, games, and other forms of relaxation have not been so often or so felicitously treated. Izaak Walton's *The Compleat Angler* (1953), however, enjoys the status of a minor classic, and the best of the modern Dutch essayists, Johan Huizinga (1872–1945), has reflected with acuteness on *Homo ludens*, or man at play. A Frenchman, Jean Prévost (1901–44), who was to die as a hero of the Resistance to the German occupation of France during World War II, opened his career as an essayist with precise and arresting analyses of the *Plaisirs des sports* (1925). But there are surprisingly few very significant works, except in chapters of novels or in short stories, on the joys of hunting, bullfighting, swimming, or even, since Anthelme Brillat-Savarin's overpraised essay, *Physiologie du goût* (1825; "The Physiology of Taste") on gourmet enjoyment of the table.

Serious speculations, on the other hand, have tended to overburden the modern essay, especially in German and in French, and to weigh it with philosophy almost as pedantic as that of academic treatises, though not as rigorous. The several volumes of Jean-Paul Sartre's *Situations*, published from 1947 on, constitute the most weighty and, in the first two volumes in particular, the most original body of essay writing of the middle of the 20th century. Albert Camus' *Mythe de Sisyphe* (1942; *Myth of Sisyphus*) and his subsequent *Homme révolté* (1951; *The Rebel*) consist of grave, but inconsistent and often unconvincing, essays loosely linked together. Émile Chartier (1868–1951), under the pseudonym Alain, exercised a lasting influence over the young through the disjointed, urbane, and occasionally provoking reflections scattered through volume after volume of his essays, entitled *Propos*.

Apart from philosophical speculation, which most readers prefer in limited quantities, the favourite theme of many modern essays has been speculation on the character of nations. It is indeed difficult to generalize on the national temper of a nation or on the characteristics of a

given culture. The authors who have done it—Emerson in his essay on *English Traits* (1856), Hippolyte Taine in his studies of the English people, Alexis de Tocqueville in his *Democracy in America* (1835, 1840)—blended undeniable conclusions with controversial assertions. Rather than systematic studies, desultory essays that weave anecdotes, intuitions, and personal remarks, ever open to challenge, have proved more effective in attempting to delineate cultures. In the 20th century, the masters of this form of writing have been among the most able in the art of essay writing: Salvador de Madariaga in Spanish, Hermann Keyserling in German, and Elie Faure in French. Some nations are much more prone than others to self-scrutiny. Several of the finest Spanish essayists were vexed by questions of what it meant to be a Spaniard, especially after the end of the 19th century when Spain was compelled to put an end to its empire. Angel Ganivet in his essay on *Idearium español* (1897; *Spain: An Interpretation*), Ortega y Gasset in *España invertebrada* (1922; *Invertebrate Spain*), and Miguel de Unamuno in almost every one of his prose essays dealt with this subject. A Spanish-born essayist, George Santayana (1863–1952), was one of the most accomplished masters of written English prose; because of his cosmopolitan culture and the subtlety of his insights, he was one of the most perceptive analysts of the English and of the American character.

Laments on the decline of the essay in the 20th century have been numerous since the 1940s, when articles in most journals tended to become shorter and to strive for more immediate effect. As a result, the general reader grew accustomed to being attacked rather than seduced. Still, the 20th century could boast of the critical essays of Virginia Woolf in England, of Edmund Wilson in America, and of Albert Thibaudet and Charles du Bos in France, all of whom maintained the high standards of excellence set by their predecessors of the previous century. It is regrettable that, in the language in which the best modern essays have been written, English, the term "essay" should also have acquired the connotation of a schoolboy's attempts at elementary composition. For the essay requires vast and varied information, yet without pedantry or excessive specialization. It must give the impression of having been composed spontaneously, with relish and zest. It should communicate an experience or depict a personality with an air of dilettantism, and of love of composition, and it should make accessible to the reader knowledge and reflection and the delight of watching a fine mind at work. The essayist should possess the virtues that one of the most influential English essayists, Matthew Arnold, praised in *Culture and Anarchy* (1869): "a passion . . . to divest knowledge of all that was harsh, uncouth, difficult, abstract, professional, exclusive; to humanize it."

HISTORY

Among the ancient Greeks and Romans, history was the branch of literature in which the most expert and the most enduring prose was written. It only recovered its supreme rank in nonfictional prose in the 18th century. Earlier, however, at the beginning of the 16th century, in Florence, Italy, Niccolò Machiavelli and Francesco Guicciardini prepared the way for history to become great literature by marrying it to the nascent science of politics and by enlarging its scope to include elements of the philosophy of history. In the 18th century Voltaire, tersely and corrosively, and Edward Gibbon, with more dignity, established history again as one of the great literary arts. In the 19th century, their lessons were taken to heart, as writers and readers realized that, in Thomas Carlyle's words, "every nation's true Bible is its history." In some nations, historians, together with epic and political poets, instilled into the people a will to recover the national consciousness that had been stifled or obliterated. Macaulay's ambition, to see history replace the latest novel on a lady's dressing table, was endorsed as an eminently reasonable and beneficent ambition by scholars throughout the 19th century. After an eclipse during the first half of the 20th century, when erudition and distrust of elaborate style prevailed, the poetry of history has again been praised by the most scrupulous practitioners of that discipline.

Essays for
entertainment

Political
essays

Inter-
pretive
history

Poetry, in that context, does not mean fiction or unfaithfulness to facts, or a mere prettification, which would be tantamount to falsification; rather, it is the recognition that, as G.M. Trevelyan, Regius professor at Cambridge University, proclaimed "The appeal of history to us all is in the last analysis poetic." Few historians today would wholly agree with the once sacrosanct formula of Leopold von Ranke (1795–1886) that their task is to record the past as it really took place. They well know that, for modern history, facts are so plentiful and so very diverse that they are only meaningful insofar as the historian selects from them, places them in a certain order, and interprets them. Since World War II, as history drew increasingly on sociology, anthropology, political and philosophical speculation, and psychoanalysis, the conviction that objectivity could be maintained by a scholar dealing with the past came to be questioned and in large measure renounced. The Italian philosopher Benedetto Croce's laconic warning that all history is contemporary history (*i.e.*, bound to the historian's time and place, hence likely to be replaced by another one after a generation) has come to be generally accepted. Nietzsche, who had sharply questioned the historical methods of his German countrymen in the 1870s, stressed the need to relate history to the present and to present it in a living and beautiful form, if it is to serve the forces of life. "You can only explain the past," he said, "by what is highest in the present."

In Germany, Italy, Spain, England, America, and most of all in France, where the vogue of sheer, and often indigestible, erudition was never wholeheartedly adopted, more literary talent may have gone into historical writing than into the novel or the short story. Many reasons account for the brilliance, and the impact, of this branch of nonfictional prose. Modern man has a powerful interest in origins—of civilization, of Christianity, of the world initiated by the Renaissance, or the French Revolution, or the rise of the masses. History invites an explanation of what is in terms of its genesis, not statically but in the process of becoming. The breadth of men's curiosity has expanded significantly since the 18th century, when belief in the absolutes of religious faith tended to be supplanted by greater concern for the relative world in which men live, move, and exist. A primary factor in the increasing importance of history is the bewilderment concerning the revolutions that occurred in or threatened so many countries in the latter part of the 20th century. As fiction, philosophy, and the exact sciences failed to provide a plausible explanation, many anguished readers turned to the record of brutal change in earlier periods. The historians who addressed themselves to those immense subjects, with their myriad ramifications, often composed monumental works of a syncretical character, such as those of Arnold Toynbee or Henri Pirenne, but they also cultivated the essay. Sometimes these essays appeared as short and pregnant volumes of reflections, such as Isaiah Berlin's *Historical Inevitability* (1954), sometimes in collections of articles that first appeared in magazines.

DOCTRINAL, PHILOSOPHICAL, AND RELIGIOUS PROSE

The question of how much of doctrinal writing, dealing with faith, ethics, and philosophy, can be called literature can only be answered subjectively by each reader, judging each case on its own merits. There have been philosophers who felt in no way flattered to be included among what they considered unthinking men of letters. The prejudice lingers in some quarters that profundity and clarity are mutually exclusive and that philosophy and social sciences therefore are beyond the reach of the layman. On the other hand, many writers, while often profound and fastidiously rigorous in their thought, such as Paul Valéry, have vehemently objected to being called philosophers. Nonetheless, a vast number of philosophical works owe their influence and perhaps their greatness to their literary merits.

In periods when philosophical speculation became very abstruse, as in Germany in the 19th century, men of letters often acted as intermediaries between the highly esoteric thinkers and the public. Much of the impact of the erudite 19th-century German philosopher Georg Wilhelm

Friederich Hegel was due to the more easily approachable writings of those who took issue with him, such as the Existentialist thinker Søren Kierkegaard, or to those who reinterpreted him, such as Karl Marx. Similarly, the thoughts of 20th-century German phenomenologist Edmund Husserl achieved wider circulation by receiving more literary expression in the writings of Jean-Paul Sartre. In modern Europe, the men of letters of Germany were long the most deeply imbued with abstract philosophy. After World War II, however, French writers appeared to take on a zest for abstract speculation, for turgid prose, and for the coining of abstruse terms. Much of French literature in the years after the war has been characterized as "literature as philosophy."

A very few philosophers have reached greatness by evolving a coherent, comprehensive system, ambitiously claiming to account for the world and man. Such harmonious constructions by the greatest philosophers, such as Descartes and Spinoza, might be compared to epic poems in sometimes embracing more than there actually appears to be between heaven and earth. These philosophical systems were conceived by powerful imaginative thinkers whose creative abilities were not primarily of an aesthetic order. The ability and the ambition to produce such systems has appeared in very few countries or cultures. The Slavic, the Spanish, and Spanish-American cultures have been richer in thinkers than in philosophers; that is, in men who reflected on the problems of their own country, who attempted to evolve a philosophy from history, or who applied a broad view to moral or political questions, rather than in men who constructed abstract philosophical systems.

More and more in the 20th century, the sciences that are called in some countries "social" and in others "humane" have replaced the all-encompassing philosophical systems of past ages. In Spain, Miguel de Unamuno (1864–1936) and José Ortega y Gasset (1883–1955) marked the thought and the sensibility of Spanish-speaking peoples far more than systematic philosophers might have done. Their writing, which disdains impeccable logic, is no less thought-provoking for being instinct with passion and with arresting literary effects.

In Russia, the doctrinal writers whose thought was most influential and often most profound were also those whose prose was most brilliant. They generally centred their speculations on two Russian preoccupations: the revival of Christian thought and charity in the Orthodox faith; and the relationship of Russia to Western Europe, branded by the Slavophiles as alien and degenerate. The consistency of ancient Greek and later Western thinkers, from Aristotle through Descartes, was of scant concern to them, but in the vitality of their style, some of these Russian theorists were masters, whose turbulent, paradoxical ideas were taken to heart by novelists, poets, and statesmen. Among these masters, Aleksandr Herzen (1812–70) combined romantic ardour and positivism, formulating a passionately Russian type of socialism; he left his mark in autobiography, political letters, fiction, and chiefly philosophy of history in *From the Other Shore* (1851). Nikolay Danilevsky (1822–85), a scientist who turned to philosophy, attempted to convince his compatriots that the manifest destiny of their country was to offer a purer and fresher ideology in lieu of that of the decadent West. V.V. Rozanov (1856–1919) was an apocalyptic prophet preaching an unusual interpretation of Christian religion; a number of his intuitions and passionate assertions are found in the novel *The Possessed* (1871–72), by Fyodor Dostoyevsky, whose own nonfictional prose is of considerable quality and conviction. The strangest and most contradictory, but also the most brilliant prose writer, among those thinkers who were torn between East and West, between a jealous Orthodox faith and the attraction of Catholic Rome, was Vladimir Solovyov (1853–1900). He blended the most personal type of visionary mysticism with an incisive humour in a manner reminiscent of Kierkegaard. His philosophical essay-dialogue-treatise, *Three Conversations on War, Progress and the End of Human History* (1900), is representative of the nonfictional Russian prose that, while not widely known outside Russia, is as revealing as the Russian novel

Russian
philosophical
essays

of the permanent contradictions and aspirations of the Slavic character.

The role of nonfictional prose in the American literature of ideas is significant, as can be seen in several of Emerson's philosophical essays and addresses; in Walt Whitman's *Democratic Vistas* (1871); in William James's pleasantly written essays on religious experience and on sundry psychological and ethical topics; in George Santayana's dexterous and seductive developments on beauty, on nature, on poets, on the genteel tradition, all envisaged with ironical sympathy. Irving Babbitt (1865–1933), Thorstein Veblen (1857–1929), and Lewis Mumford are among the many American writers who, in the 20th century, maintained the tradition of writing on abstract or moral themes with clarity and elegant simplicity. Earlier, Thomas Jefferson and Benjamin Franklin had expressed their lay philosophy in a manner they wished to be widely accessible.

In France the tradition *haute vulgarisation*—"higher vulgarization" or popularization—never died and was seldom slighted by the specialists. There, and to a slightly lesser extent in Britain, much of the most valuable writing in prose was an elucidation of the view of life underlying the creations of eminent men in many fields. Such doctrinal writing, expounding innermost convictions and sometimes representing a diversion from more intensive pursuits, constitutes a by no means negligible portion of the writings of the philosopher Bertrand Russell, of the poet William Butler Yeats, and others. The novelist or the poet may well use nonfictional prose to purge his own anger, to give vent to his vituperation against his confrères, and to relieve his imagination of all the ideological burden that might otherwise encumber it. D.H. Lawrence preserved the purity of his storyteller's art by expressing elsewhere his animadversions against Thomas Hardy or Sigmund Freud. Albert Camus stripped his fiction and short stories of the ideological musings found in his philosophical volumes. Marcel Proust succeeded in incorporating many abstract discussions of the value of art, love, and friendship in his very original and loose type of fiction; but his great work, *À la recherche du temps perdu* might well have gained even more from the excision of those dissertations and the writing of more volumes like his *Chroniques* (1927) or *Contre Sainte-Beuve* (1954). The masters of nonfictional prose in French in the 20th century have been those thinkers who were also superb stylists and who deemed it a function of philosophy to understand the aesthetic phenomenon: Henri Bergson (1859–1941), Paul Valéry (1871–1945), and Gaston Bachelard (1884–1962). No more poetical advocate of reverie has arisen in the 20th century than *La Poétique de la rêverie* (1960; *The Poetics of Reverie*) and the posthumous collection of essays, *Le Droit de rêver* (1970; "The Right to Dream"), by Bachelard, who was also a philosopher of science. A major influence on him, as on several earlier poets endowed with profound intellect, such as Baudelaire and Valéry, was Edgar Allan Poe, the impact of whose essays on poetics, on cosmology, and on dreams and reveries has been immense and beneficent. More than a century after his death, many of Poe's American compatriots have conceded that the storyteller and the poet in Poe counted for less, as his European admirers had divined, than the writer of critical and doctrinal prose rich in dazzling intuitions.

Although lectures, articles, and other prosaic admonitions have tended to take their place, sermons, funeral orations, allegories, and the visions of eternal punishment brandished by theologians constitute some of the most unforgettable prose. This form of nonfictional prose literature dates from before the Christian Era; Jewish thought and style were molded by commentaries on the Old Testament and compilations of the wisdom of the sages. Later, and more nearly literary, works of this nature include Sebastian Brant's didactic, poetical, and satirical *Narrenschiff* (1494; *Ship of Fools*), and the mystic writings of Jakob Böhme (1575–1624) in Germany, the moving sermons of Jón Vídalín (1666–1720) in Iceland. In England, Richard Baxter (1615–91) and John Bunyan (1628–88) were among the most eloquent of the 17th-century Puritans who composed doctrinal works of literary merit; along

with the epic poet John Milton (1608–74), whose prose works hardly count for less than his poetry, they exercised a powerful influence on the English language through their doctrinal prose. Their contemporary, the Anglican Jeremy Taylor (1613–67), wrote the most sustained and dignified prose of an age that, on the continent, would be called Baroque. A little later, in northern Europe, the Norwegian Ludvig Holberg (1684–1754), who spent most of his life in Denmark and became best known as a comic writer, also advised his contemporaries how to live morally in his *Ethical Thoughts* and other didactic treatises. The Swede Emanuel Swedenborg (1688–1772), less gifted as a writer but far more original in his blend of mysticism and science, outshone all previous Scandinavians in impressing the imagination of other Europeans. No less influential, Søren Kierkegaard (1813–55), because of his stimulating ambiguities, his bold treatment of traditional theology and philosophy, and his extraordinary ability to write vivid, biting, and provoking prose, was, a century after his death, one of the most potent forces in the literature and thought of Western civilization.

Many 20th-century readers experience a feeling of remoteness in this kind of doctrinal writing, which stems in part from a lack of vital interest in the beliefs it embodies and from a coolness toward religious dogmatism or fanaticism. Intolerance has shifted from religion to the domain of politics. But the contemporary estrangement from that rich literary heritage is due also to a distrust of high-flown eloquence. Cotton Mather's *Essays to do good* (1710) has few readers in present-day New England, despite that region's Puritan tradition, and Jonathan Edwards (1703–58), a writer of great spiritual warmth and imaginative style who was the first of the great prose writers of America, is admired today chiefly by specialists.

A less sonorous style, one that does not ring so monotonously ornate to the reader's ears, is now preferred. In Spain, Antonio de Guevara (c. 1480–1545), a preacher who was at his best in his familiar and satirical moments, and St. Teresa of Avila (1515–82), in her records of her mystical ecstasies, have withstood the changing tides of taste. The French also succeeded in maintaining their appreciation of their two greatest religious writers, Pascal and Bossuet, at the very top of the nonfictional prose writers; both are still revered and occasionally imitated. Pascal took over traditional theology and treated it as literature; his unfinished *Pensées* have exercised far more influence than the rationalism of the greatest French philosophers on the sensibilities of the French. Bossuet's orations reveal the magnificent but refrigerating decorum that seems inseparable from eulogies of the dead—a genre that precludes full sincerity and cultivates tremulous emotion to a dangerous degree. Bossuet's sermons and treatises, however, include masterpieces of simple, terse, direct oratory, which show him as the majestic defender of the unity of faith, of absolutism, and of tradition. His was the last significant endeavour in the 17th century to arrest the flow of relativism and of rebellious individualism, which had engulfed Western civilization with the Renaissance, the Reformation, and Humanism. The two most brilliant writers of religious prose in France in the 20th century were Pierre Teilhard de Chardin (1881–1955), a poetical writer with a luxury of images, and Simone Weil (1909–43), more terse and restrained; they steered a middle course between dogmatism and humility in luring the lay reader to their ardent expressions of conviction.

POLITICAL, POLEMICAL, AND SCIENTIFIC PROSE

In the 20th century, political, economic, and social thought has attempted to reach scientific precision through the use of quantitative data, processing machines, and mathematical formulas. Through such means, other disciplines eventually were elevated to the status of sciences. Literature lost a great deal as a result of this scientific urge, and political and economic thought may have lost even more; for example, the ability to be understood, and perhaps applied, by men of affairs and leaders of nations. The result has been that momentous decisions may be made independent of political theory, which is more often called upon to explain them afterward. Albert Einstein

remarked that politics is much more baffling and difficult than physics and that consequences of errors in politics are likely to make far more difference to the world than the miscalculations of science. Politics is often defined as the art of the possible; it is also an art of improvisation, since the fleeting occasions must be grasped when they appear, and risks must be taken without a full array of scientific data. Like military action, however, political action can be studied in historical writings and in the literary testimonials of men who ran the affairs of their country. Thucydides, Cicero, Caesar, Milton, Burke, Napoleon, and Jefferson were such men of action who were also endowed with uncommon literary gifts. In varying degrees, Benjamin Disraeli, Winston Churchill, Woodrow Wilson, Clemenceau, Lenin, and de Gaulle owed some of their insight and effectiveness to their literary efforts.

Authors, however, are by no means infallible in dealing with the unpredictable course of political life. Interpreting and channelling public opinion proved insuperably difficult, for example, to Alphonse de Lamartine in the revolutionary period of 1848–49 in France, to the bookish Aleksandr Kerensky during the 1917 Revolution in Russia, and to a number of brilliant writers who attempted to guide the Spanish Republic in the 1930s. Crowds often can be moved more readily by vapid, repetitious, or inflammatory speeches than by profound or wise counsel. Abraham Lincoln's Gettysburg Address, Churchill's speeches during Britain's "finest hour" early in World War II, and de Gaulle's lofty eloquence regarding the crises of three decades in France were admired less when they were delivered than afterward. As they are collected, studied, and engraved in the mental makeup of millions of future citizens, such speeches have an effectiveness second to no other form of nonfictional prose. Novels may exercise immense influence through the acute social criticism they embody, but their impact upon the sensibility and the behaviour of their readers is probably less than that of political prose.

Although the Spanish language cannot boast of any political thinker comparable to Plato, Machiavelli, or Rousseau, it may boast a large number of fine writers on political topics. Generally, these writers reveal a restrained and terse style, like the poets of Spain, the Latin country least addicted to inflation of language. Garcilaso de la Vega (1539–1616), the son of an Inca mother, wrote with courage and talent of the Peruvians and other cultures of the New World cruelly wrecked by their Catholic conquerors. The Argentinian Domingo Faustino Sarmiento (1811–88) fought in battle and with his pen against his country's dictator and left a masterpiece of social insight, written with rare effectiveness, *Facundo* (1845). Miguel Ángel Asturias (1899–1974), from Guatemala, scathingly depicted the evils of dictatorship in Central America. Like many others in South America, where versatility is not uncommon, Francisco de Miranda (1750–1816) of Venezuela was both a political writer and a statesman.

Italy, after Machiavelli, failed to produce political writers of very great eminence, even during its liberation and unification in the 19th century. The universal thinker Benedetto Croce (1866–1952), however, had the courage to publish during the Fascist era the most impassioned defense of liberty in volumes such as *La storia come pensiero e come azione* (1938; *History as the Story of Liberty*). Another Italian—but from another political direction—Antonio Gramsci (1891–1937), one of the most intelligent exponents of Communism in western Europe, was aware of the vital significance of literary form to spread political ideas. He bitterly deplored the lack of a popular literature in his country that reflected the morality and sentiment of the people.

In France political speculation was more comprehensive: few political theoreticians have proved as influential as the philosophers of the Enlightenment, especially Montesquieu and Rousseau. It was the good fortune of the French that during their Revolution at the end of the 18th century and throughout the 19th century, its keenest political minds were also writers of admirable prose. Tocqueville's observations became a sacred text for many a student of America and of pre-Revolutionary France.

Since the French seldom give ideas serious consideration unless they are well expressed, however, it was a misfortune that most political speculation after the Napoleonic age was written by gifted, often brilliant, conservatives, such as Joseph de Maistre, Auguste Comte, Frédéric Le Play, Renan, Taine, and Charles Maurras. Those advocating a socialistic view, such as Jean Jaurès and the more elegant and genteel Léon Blum, failed to express their theories in classic prose. The level of political comment in the magazines and newspapers in France is consistently high, but the writers tend to be either too clear-sighted or too arrogant to grant their statesmen a chance to act. "Fair play" is an untranslatable phrase in French, and politics in France, unlike some other countries, is never regarded as a game or sport. Rather, it is a passionate affair of the heart and intellect, conducted in a mood of intransigence. The English essayist Walter Bagehot (1826–77), observing the French at the time of the 1851 coup d'état, commented wryly that "the most essential quality for a free people, whose liberty is to be progressive, permanent and on a large scale, is much stupidity . . . Stupidity is nature's favorite resource for preserving steadiness of conduct and consistency of opinion."

English and American political works, from the 17th century on, excel all others; they constitute the richest form of nonfictional prose in the English language. John Milton's *Areopagitica* (1644) and his other political pamphlets are monuments of political prose that survive to this day as classics. Edmund Burke's *Letter to a Noble Lord* (1796) was praised a century and a half after its composition as the greatest piece of invective in the English language. William Godwin's *Political Justice* (1793) does not compare in the majesty of its prose to those supreme models, but it did inflame Shelley and other men of letters of the time. Walter Bagehot wrote equally well on literature, politics, and economics, and *The Economist*, which he edited, was the best-written weekly of its kind in any language. John Stuart Mill and Thomas Carlyle also helped to maintain the tradition of political and social thought expressed as literature through the 19th century.

Polemical prose significantly declined in the modern era. Few moderns express the rage for invective seen in the verse of satirists such as the ancient Roman Juvenal or Alexander Pope in 17th-century England or even in the writings of Christian disputants such as Martin Luther. Voltaire rejoiced in flaying not only his enemies but also some, such as Montesquieu and Rousseau, who were fundamentally in agreement with him in the fight against the religion of his age. Literary polemics of a high order were employed against the cultural imperialism of the French in Gotthold Lessing's *Hamburgische Dramaturgie* (1767–69; *Hamburg Dramaturgy*). Beside these examples, the polemics of more recent periods seem tame, or else gross and venomous. Later practitioners of the literature of insult include Émile Zola, particularly in his celebrated article on the Dreyfus affair, *J'Accuse* (1898). Later writers, however, often overreach themselves; their rhetoric sounds vapid and their epigrams strained.

The rift between the two cultures, scientific and humanistic, is probably not as pronounced or final as it has been alleged to be. About the time the division was enunciated, in the mid-20th century, it was possible to point to a number of eminent scientists who were also masters of prose writing—Henri Poincaré, Jean Rostand, and Gaston Bachelard in France; Bertrand Russell and Alfred North Whitehead in England; and René Dubos and Robert Oppenheimer in the U.S. The peril for scientists who undertake to write for laymen appears to lie in a temptation to resort to florid language and to multiply pretentious metaphors and elaborate cadences in their prose. Some scientists who wrote on astronomy, on anthropology, and on geology have not altogether escaped that pitfall: Sir James Jeans, Loren Eiseley, Sir James Frazer, Teilhard de Chardin. The marriage of the "two cultures" in one mind, which was no less concerned with scientific truth than with beauty of form, was found frequently in older times; Aristotle, Hippocrates, Galileo, Newton, and Goethe all showed strong interest in both. The popularization of science reached a level of a lucid and elegant art with the

writings of Bernard de Fontenelle (1657–1757) in French, Francesco Algarotti (1712–64) in Italian, and later, with a masterpiece of scientific rigour expressed in flexible and precise prose. *Introduction à l'étude de la médecine expérimentale*, by the physiologist Claude Bernard (1813–78).

OTHER FORMS

Reportage. Journalism often takes on a polemical cast in countries in which libel laws are not stringent. Polemical journalism flourished in continental Europe when a journalist's insults could be avenged only in a duel: one of the great journalists of this heroic era of the press in France, Armand Carrel, died in such a duel with another journalist in 1836. Most journalistic literature, however, deserves none of the ill-repute that is associated with its more polemical expressions. Rather, it is a remarkably elastic form, as adaptable to sarcasm and the puncturing of illusions as to reflection, subtle persuasion, and infectious geniality. Among the eminent writers who explored its possibilities in the 18th century, Joseph Addison and Richard Steele offered models of polished English prose in the journals *The Tatler* and *The Spectator*, and Jonathan Swift and Oliver Goldsmith also used it effectively in England. In France Voltaire, the novelist Abbé Prévost, and the dramatist Pierre-Carlet de Marivaux all found effective use for the form. By the 19th century, most eminent men of letters attempted to broaden their audiences by means of articles and essays in the press, and in the 20th century, the influence of journalism pervaded the most important works of some authors. Some of the works of G.B. Shaw and H.G. Wells, for example, were reminiscent of journalism in the manner in which they sought topical controversy and challenged social and political prejudices. Many of the finest essays of Virginia Woolf, John Middleton Murry, and Aldous Huxley represented British literary journalism at its most intelligent level. In America, the more heterogeneous public to which authors must address themselves and, later, the competition of the audiovisual media, were not propitious to the flowering of literary journalism of that type. In a more ephemeral genre, that of political reflections couched in clear, pungent style, Walter Lippmann composed models of commentaries on politics and ethics.

Ego in
journalistic
style

The more self-centred and passionate writers seldom succeeded in journalistic prose as well as those who could forget their ego and adapt their style to a public that wanted to be entertained, moved, or convinced, perhaps, but whose attention span extended no further than the 15 minutes of a train ride or of a hurried breakfast. In France, Proust dreamt for years of appearing as a journalist on the first column of the journal *Le Figaro*. But he and his contemporaries Gide, Claudel, and Valéry, and, later, the imperious and nervous André Malraux, did not conform to the limitations of the newspaper article. On the other hand, Colette, Paul Morand, and François Mauriac proved conspicuously successful in writing the brief, gripping, taut article dear to readers of many of the better continental dailies and weeklies.

The insidious appeal of journalistic writing to thinkers, novelists, and poets is similar to the siren charm of conversation for the author who enjoys talking brilliantly at dinner parties. As Oscar Wilde ruefully remarked, conversationalists and journalists, intent on reporting on the ephemeral, pour whatever genius is theirs into their lives, and only their talent into their works.

Aphorisms and sketches. Authors of maxims and aphorisms, on the contrary, strive for the brevity of inscriptions on medals and public buildings and for a diamond-like resistance to the devastation of time upon diffuse and padded writing. This form is periodically revived. In modern letters, in the latter half of the 20th century, a condensed and enigmatic sort of prose was preferred to poetry by several poets, who invested their sensations, their illuminations, or their reflections with the mystery and éclat of aphorisms. Among the French, who have always favoured the maxim for philosophical, psychological, and ethical advice, a great poet, René Char, came to be more and more fascinated by that epigrammatic form, harking

back to the ancient Greek philosopher whom he admired most, Heraclitus. Char found in the aphorism a means of "pulverizing language" and of allowing the isolated words or groups of words, freed from rhetoric and from the exigencies of clarity, to emerge like rocks from a sunken archipelago. Other French prose writers, including Camus, Char's warmest admirer, and Malraux, likewise scattered through their prose works striking aphorisms that summed up the sense of a situation or the experience of a lifetime. French novels, from the 18th century through the 20th, reflect the influence of the unforgettable maxims coined by the 17th-century moralists Pascal, La Rochefoucauld, and La Bruyère. The novelist could never long resist the seduction of brevity, the challenge of condensing wisdom into a neat, usually bitter, formula, which usually suggested to the reader not to expect overmuch from life and to take revenge upon its little ironies by denouncing it in advance.

Maxims and other pointed and epigrammatic phrases of the sort the ancient Romans called *sententiae* can become too sophisticated or can too obviously strive for effect. This form of expression reached its point of perfection, balancing profundity and solidity of content with pointedness of form, with the moralists of the 17th and 18th centuries in France, whom Nietzsche ranked above all other writers. They included Pascal and La Rochefoucauld and, later, Sébastien Chamfort (1740/41–94), a satirical pessimist often quoted by Schopenhauer and Joseph Joubert (1754–1824). This form, even more than poetry, represents the most economical means of communicating long experience and for imparting moral advice. In a very few words, or at most a few lines, an aphorism may enclose enough matter for the plot of a novel. It may trounce the prejudices of snobbery more vigorously than a long, meandering novel of manners. The greatest of the 19th-century poets, Goethe, Novalis, Leopardi, Vigny, and Baudelaire, as well as painters such as Delacroix, Cézanne, Degas, and later Braque, cherished the epigrammatic, incisive form of expression. One of the advantages of the aphorism or *pensée* is that it can easily produce an impression of depth when it may be only a commonplace pungently expressed. Another is that it allows several approaches to a subject by the skilled prose writer. If he is of a fiery temperament, prone to enthusiasms and lashing out in wrath against what he deems to be false, he can, like Nietzsche, embrace contradictions and sponsor opposed attitudes. If Epictetus, Pascal, and Nietzsche had expressed their reflections consistently and systematically, their works would probably be forgotten. Nonetheless, as Pascal shrewdly remarked, the aphoristic prose style is, of all the manners of writing, the one that engraves itself most lastingly in the memories of men.

That form, in verse and in prose, probably constitutes the most widespread form of literature. It is found in many nations that long lived without fiction, epics, or even popular poetry. It is found in ancient sayings that interlard the speeches of the 20th-century leaders both of the U.S.S.R. and of China; in the book of Proverbs of the Bible; in the Qur'an; in the Afrikaans language of South Africa in the 20th-century writings of J. Langenhoven. Proverbs, maxims, riddles, and even conundrums make up a large part of black African folklore. The animal tales of these people also provide lessons in the form of aphorisms that are neither as platitudinous nor as didactic as Aesop's fables.

Portraits and sketches are a form of literature that thrives in cultures in which the court, the salon, or the café plays an important role. The few examples left by the ancient Greeks, such as by Theophrastes, pale beside the vivid portraits of real individuals drawn by the ancient Roman historian Tacitus and by the impassioned orator Cicero. In the Classical age of 17th-century France, the character sketch was cultivated in the salons and reached its summit with La Bruyère. That form of writing, however, suffered from an air of artificiality and of virtuosity. It lacked the ebullience and the imagination in suggesting telltale traits that characterize the portraits of the duc de Saint-Simon (1675–1755). Collections of sketches and characters, however, tend to strike the reader as condescending

Portraits
and
character
sketches

and ungenerous insofar as the writer exempts himself of the foibles he ridicules in others.

The humorous article or essay, on the other hand, is a blend of sympathy and gentle pity with irony, a form of criticism that gently mocks not only others but the mocker himself. Humour strikes deep roots in the sensibility of a people, and each nation tends to feel that its own brand of humour is the only authentic one. Its varieties of humorous writing are endless, and few rules can ever be formulated on them. Humorous literature on the highest literary level includes that of Cervantes in Spain, of Sterne, Lamb, and Thackeray in England, of Jean Paul in Germany, and of Rabelais, Montaigne, and Voltaire in France. Romantic authors have, as a rule, been too self-centred and too passionate to acquire the distance from their own selves that is essential to humour. In the 20th century, some of the most original examples of what has been called the "inner-directed smile" are in the works of the Argentine Jorge Luis Borges and by one of the writers he admires most, the English essayist G.K. Chesterton (1874–1936). In both writers, and in other virtuosos of the intellectual fantasy, there is a persistent refusal to regard themselves as being great, though greatness seems to be within their reach. The humorist will not take himself seriously. Chesterton hides the depth of his religious convictions, while Borges facetiously presents his prodigious erudition and indulges in overelaborate and flowery prose. Borges likes to put on and take off masks, to play with labyrinths and mirrors, but always with a smile. By sketching what appear to be fanciful portraits rather than overtly fictional stories, he creates a half-imaginary character whose presence haunts us in all his writings—that of the author himself.

Dialogues. The dialogue form has long been used as a vehicle for the expression of ideas. It is especially cherished by authors eager to eschew the forbidding tone of formality that often accompanies the expression of serious thought. The writer of a dialogue does not directly address his public, but instead revels in the multiple facets of ideas. By playing this dialectical game he can appear to present contrary views as their respective proponents might and then expose the errors of those he opposes, leading the readers to accept his own conclusions. The advantages of the dialogue are clear: ideas that might have remained abstruse and abstract become concrete and alive. They assume dramatic force. A constant element in the dialogue is irony; etymologically, the term derives from a form of interrogation in which the answer is known beforehand by the questioner. The earliest models of the genre, by the ancient Greeks Plato and Lucian, have never been excelled. Sophistry is another element of the dialogue. In Plato and in the dialogues of Pascal's *Provinciales* (1656–57: "Provincial Letters"), the protagonist plays with the naiveté of his opponents, who always end by surrendering. The writer of a dialogue cannot affect the same casual and self-indulgent attitude as the author of a personal essay since the characters and their statements must be plausible. Nor can he pursue an argument consistently, as he might in a critical, historical, or philosophical essay. Something must persist in the dialogue of the spontaneity and the versatility of an actual conversation among witty and thoughtful people.

There was much seriousness and occasionally some pedantry in early dialogues in several literatures. The dialogues of Bardesanes (154–222) in Syriac, rendered into English as *On Fate*, are on the subject of the laws of the country. A hundred years earlier, Lucian, who was also Syrian, proved himself a master of flowing and ironical Greek prose in his satirical dialogues. The Italian Renaissance writer Pietro Aretino (1492–1556) proved himself the equal of Lucian in verve in his *Dialogues* using the same mold and the same title as Lucian. Others who used the dialogue form included Castiglione and Pietro Bembo (1470–1547) in Italy; and in Spain Juan Luis Vives (1492–1540), León Hebreo (1460–c. 1521), and Juan de Valdés (c. 1500–41), who treated questions of faith and of languages in dialogues. The genre flourished in the 18th century: Lessing, Diderot, and the Irish philosopher George Berkeley. Diderot's works largely consist of sprightly, rambling, and provocative discussions between

the various aspects of his own remarkable mentality. Bold conjectures, determined onslaughts on prejudices, insights into physiology and biology, and erotic fantasies all enter into his dialogues. In the 19th century a number of complex literary personalities, who were capable of accepting the most diverse, and even conflicting points of view, such as Renan and Valéry, had a predilection for the dialogue. Among the devices used by authors of dialogue—many of whom lacked the sustained inventiveness required by fiction—was to attribute their words to the illustrious dead. The French prelate Fénelon, for example, composed *Dialogues des morts* (1700–18), and so did many others, including the most felicitous master of that prose form, the English poet Walter Savage Landor, in his *Imaginary Conversations* (1824) and *Pentameron* (1837).

Travel and epistolary literature. The literature of travel has declined in quality in the age when travel has become most common—the present. In this nonfictional prose form, the traveller himself has always counted for more than the places he visited, and in the past, he tended to be an adventurer or a connoisseur of art, of landscapes, or of strange customs who was also, occasionally, a writer of merit. The few travel books by ancient Greek geographers, such as Strabo and Pausanias of the 1st and 2nd centuries AD, are valuable as a storehouse of remarks on ancient people, places, and creeds. Travel writing of some literary significance appears in the late-13th-century writings of Marco Polo. Works of a similar vein appeared in the 17th century in the observations of Persia two French Huguenots, Jean-Baptiste Tavernier and Jean Chardin, whose writings were lauded by Goethe. Many books of documentary value were later written by English gentlemen on their grand tour of the Continent. The 18th-century Italian egotist Casanova and his more reliable and sharper compatriot Giuseppe Baretti (1719–89) also produced significant travel writings.

The form comprises many of the finest writings in prose during the Romantic age. Not only were the Romantics more alive to picturesqueness and quaintness but also they were in love with nature. They were eager to study local colours and climates and to depict them in the settings for their imaginative stories. Also, travel gave the Romantic writer the illusion of flight from his wearied self. The leisurely record of Goethe's journey to Italy in 1786–88 counts more readers than most of his novels. *Pismo russkogu puteshestvennika* (1791–92; Eng. trans., *Letters of a Russian Traveler, 1789–1790*, 1957) by Nikolay Karamzin is one of the earliest documents in the development of Russian Romanticism. Ivan Goncharov (1812–91), the Russian novelist who stubbornly limited his fiction to his own geographical province, recorded in *Frigate Pallas* his experience of a tour around the world. Nowhere else in the whole range of literature is there anything comparable to *Peterburg* (1913–14), by a virtuoso of poetic style, Andrey Bely; it is a travel fantasy within a city that is both real and transfigured into a myth. Neither James Joyce's Dublin nor Balzac's Paris is as vividly recreated as the former Russian capital in Bely's book. Other travel writers of note include the multinational Lafcadio Hearn (1850–1904), who interpreted Japan with sensitivity and insight. Earlier, two other Westerners wrote on Asia, the English historian Alexander W. Kinglake (1809–91), in *Eothen* (1844), and, more incisively, the French diplomat Joseph-Arthur, comte de Gobineau (1816–82); both blended a sense of the picturesqueness of the East with shrewdness in the interpretation of the people. One of the most thoughtful and, in spite of the author's excessive self-assurance, most profound books on Asia is *Das Reisetagebuch eines Philosophen* (1919; *Travel Diary of a Philosopher*), by the German thinker Hermann Keyserling (1880–1946). With an insatiable interest in countries, Keyserling also interpreted the soul of South America and, less perceptively, analyzed the whole spectrum of European nations. Among the thousands of travel books on Italy, there are a few masterpieces of rapturous or humorous prose: in English, the writings of D.H. Lawrence on Sardinia, on Etruscan Italy, and on the Italian character are more lucid and less strained than other of his prose cogitations. Venice, "man's most beautiful artifact," as Bernard Beren-

son called it, inspired Rousseau, Chateaubriand, Maurice Barrès, Anatole France, and hundreds of other Frenchmen to write some of their finest pages of prose. After World War I, there was a distinct yearning for new possibilities of salvation among war-ridden Europeans, dimly described in Asia, in Russia, or in America, and travel literature assumed a metaphysical and semireligious significance. The mood of the writers who expressed this urge was somewhat Byronic; they were expert at poetizing the flight from their own selves. Blaise Cendrars (1887–1961) in his novel *Emmène-moi au bout du monde* (1956: “Take Me Away to the End of the World”), epitomizes the urge to seek adventures and a rediscovery of oneself through strange travels. The very theme of travel, of the protagonist being but a traveller on this earth, has been, from Homer’s *Odyssey* onward, one of the most laden with magical, and symbolical, associations in literature. Countless authors have played moving and delicate variations on it.

The letter
as a genre

Of all the branches of nonfictional prose, none is less amenable to critical definition and categorization than letter writing. The instructions of the ancient grammarians, which were repeated a thousand times afterward in manuals purporting to teach how to write a letter, can be reduced to a few very general platitudes: be natural and appear spontaneous but not garrulous and verbose; avoid dryness and declamatory pomp; appear neither unconcerned nor effusive; express emotion without lapsing into sentimentality; avoid pedantry on the one hand and banter and levity on the other. Letters vary too much in content, however, for generalizations to be valid to all types. What is moving in a love letter might sound indiscreet in a letter of friendship; an analysis of the self may fascinate some readers, while others prefer anecdotes and scandal. La Bruyère, at the end of the 17th century, remarked that women succeed better than men in the epistolary form. It has also been claimed that a feminine sensibility can be seen in the letters of the most highly acclaimed male masters of this form, such as Voltaire, Mirabeau, Keats, and Baudelaire. Advice to practitioners of the art of letter writing usually can be expressed in the often-quoted line in Shakespeare’s *Hamlet*: “To thine own self be true.” The English biographer Lytton Strachey (1880–1932), a copious and versatile letter writer himself, wrote: “No good letter was ever written to convey information, or to please its recipient: it may achieve both those results incidentally; but its fundamental purpose is to express the personality of the writer.” There are, however, numerous and even contradictory ways of expressing that personality.

Although critics have issued endless disquisitions on the craft of fiction and other genres, they have generally remained silent on the epistolary genre, though it has sometimes been the form of prose that outlives all others. Ever since the expression of the writer’s personality became one of the implicit purposes of writing in the 18th century, the letters of such eminent authors as Diderot, Rousseau, Byron, and Flaubert have probably offered at least as much delight as any of their other writings. Impressive monuments of scholarship have been erected on the presentation of the complete letters of Thackeray, George Eliot, Swinburne, and Henry James. The literatures of France and England are notably richer in letter writing of the highest order than are the literatures of the United States and Germany. Contrary to many pessimistic predictions regarding the effect on letter writing of modern means of communication, such as the telephone, together with an apparently increasing penchant for haste, some of the richest, most revealing, and most thoughtful letters of all times have been written in the 20th century: those of the English writers Katherine Mansfield and D.H. Lawrence are paramount among them.

Personal literature. The cult of the ego (that is, a preoccupation with self-analysis) is a late development in the history of literature. There were, to be sure, men in ancient times who were absorbed in their own selves, but there is almost no autobiographical literature from ancient Greece and, in spite of Cicero and Pliny the Younger, there is little from ancient Rome. The confession, made as humble as possible and often declamatory in the exposition of the convert’s repented sins, was an outgrowth of

Christianity; masters of confessional literature were Saint Augustine, Petrarch, and the English Puritans. Autobiographical writing took a different form in the 18th century in the work of men who would have agreed with Goethe that personality is the most precious possession. After the publication of Rousseau’s *Confessions* in France in 1781, the passion for looking into one’s heart (and other organs as well) spread to other literatures of western Europe. Many a novelist thereafter kept a precise record of his cogitations, anxieties, and harrowing moments of inability to create. Poets and painters, including Delacroix, Constable, and Braque, have often done the same. There is only a very tenuous separation between fiction of this sort from nonfiction; the introspective novel in the first person singular has much in common with a diary, or a volume of personal reminiscences. In his long novel *À la recherche du temps perdu*, Proust revealed himself in three ways—as the author, as the narrator, and as the characters who are projections of his own self. An autobiography once was ordinarily written toward the end of a life, as a fond recollection or an impassioned justification of a lifetime’s deeds. More and more, it has come to be written also by men and women in their prime. The names of writers whose autobiographical writings have become classics is legion. Henry Adams (1838–1918) owes his place in American letters chiefly to his book on his education; in 20th-century English letters, Osbert and Sacheverell Sitwell, Leonard Woolf, and Stephen Spender may similarly survive in literature through autobiographical works. André Gide, always uncertain of his novelist’s vocation, felt more at ease laying bare the secret of his life in autobiographies and journals.

Although imaginative fiction has probably suffered from excesses of introspection and of analyses of the author’s own artistic pangs, knowledge of man’s inner life has been enriched by such confessions. The most profound truths on human nature, however, have been expressed not in the form of autobiography but in its transposition into fiction. Readers generally have found more truth in literature created from the possibilities of life than from the personal record of the one life that the author has lived.

In conclusion, the variety of nonfictional prose is prodigious. It can be written on almost any conceivable subject. Almost any style may be used, from casual digressions or sumptuous and sonorous sentences to sharp maxims and elliptical statements. But nonfictional prose seldom gives the reader a sense of its being inevitable, as does the best poetry or fiction. Nonfictional prose seldom can answer positively the question that Rilke and D.H. Lawrence suggest that any potential writer should ask: Would I die if I were prevented from writing? (H.M.P.)

Biographical literature

One of the oldest forms of literary expression, biographical literature seeks to recreate in words the life of a human being, that of the writer himself or of another person, drawing upon the resources, memory and all available evidences—written, oral, pictorial.

ASPECTS

Historical. Biography is sometimes regarded as a branch of history, and earlier biographical writings—such as the 15th-century *Mémoires* of the French councillor of state, Philippe de Commines, or George Cavendish’s 16th-century life of Thomas Cardinal Wolsey—have often been treated as historical material rather than as literary works in their own right. Some entries in ancient Chinese chronicles included biographical sketches; imbedded in the Roman historian Tacitus’ *Annals* is the most famous biography of the emperor Tiberius; conversely, Sir Winston Churchill’s magnificent life of his ancestor John Churchill, first duke of Marlborough, can be read as a history (written from a special point of view) of Britain and much of Europe during the War of the Spanish Succession (1701–14). Yet there is general recognition today that history and biography are quite distinct forms of literature. History usually deals in generalizations about a period of time (for example, the Renaissance), about a group of people

in time (the English colonies in North America), about an institution (monasticism during the Middle Ages). Biography focusses upon a single human being and deals in the particulars of his life.

Both biography and history, however, are concerned with the past, and it is in the hunting down, evaluating, and selection of sources that they are akin. In this sense biography can be regarded as a craft rather than an art: techniques of research and general rules for testing evidence can be learned by anyone and thus need involve comparatively little of that personal commitment associated with art.

A biographer in pursuit of an individual long dead is usually hampered by a lack of sources: it is often impossible to check or verify what written evidence there is; there are no witnesses to cross-examine. No method has yet been developed by which to overcome such problems. Each life, however, presents its own opportunities as well as specific difficulties to the biographer: the ingenuity with which he handles gaps in the record—by providing information, for example, about the age that casts light upon the subject—has much to do with the quality of his resulting work. James Boswell knew comparatively little about Dr. Johnson's earlier years; it is one of the great-nesses of his *Life of Samuel Johnson LL.D.* (1791) that he succeeded, without inventing matter or deceiving the reader, in giving the sense of a life progressively unfolding. Another masterpiece of reconstruction in the face of little evidence is A.J.A. Symons' biography of the English author and eccentric Frederick William Rolfe, *The Quest for Corvo* (1934). A further difficulty is the unreliability of most collections of papers, letters, and other memorabilia edited before the 20th century. Not only did editors feel free to omit and transpose materials, but sometimes the authors of documents revised their personal writings for the benefit of posterity, often falsifying the record and presenting their biographers with a difficult situation when the originals were no longer extant.

The biographer writing the life of a person recently dead is often faced with the opposite problem: an abundance of living witnesses and a plethora of materials, which include the subject's papers and letters, sometimes reports of telephone conversations and conferences transcribed from tape, as well as the record of interviews granted the biographer by his subject's friends and associates. Frank Friedel, for example, in creating a biography of the United States president Franklin D. Roosevelt (1882–1945), has had to wrestle with something like 40 tons of paper. But finally, when writing the life of any man, whether long or recently dead, the biographer's chief responsibility is vigorously to test the authenticity of his materials by whatever rules and techniques are open to him.

Psychological. Assembling a string of facts in chronological order does not constitute the life of a person, it only gives an outline of events. The biographer therefore seeks to elicit from his materials the motives for his subject's actions and to discover the shape of his personality. The biographer who has known his subject in life enjoys the advantage of his own direct impressions, often fortified by what the subject has himself revealed in conversations, and of his having lived in the same era (thus avoiding the pitfalls in depicting distant centuries). But on the debit side, such a biographer's view is coloured by the emotional factor almost inevitably present in a living association. Conversely, the biographer who knows his subject only from written evidence, and perhaps from the report of witnesses, lacks the insight generated by a personal relationship but can generally command a greater objectivity in his effort to probe his subject's inner life.

Biographers of the 20th century have had at their disposal the psychological theories and practice of Sigmund Freud and of his followers and rivals. The extent to which these new biographical tools for the unlocking of personality have been employed and the results of their use have varied greatly. On the one hand, some biographers have deployed upon their pages the apparatus of psychological revelation—analysis of behaviour symbols, interpretation based on the Oedipus complex, detection of Jungian archetypal patterns of behaviour, and the like. Other bio-

graphers, usually the authors of scholarly large-scale lives, have continued to ignore the psychological method; while still others, though avoiding explicit psychological analysis and terminology, have nonetheless presented aspects of their subjects' behaviours in such a way as to suggest psychological interpretations. In general, the movement, since World War I, has been toward a discreet use of the psychological method, from Katherine Anthony's *Margaret Fuller* (1920) and Joseph Wood Krutch's study of Edgar Allan Poe (1926), which enthusiastically embrace such techniques, through Erik Erikson's *Young Man Luther* (1958) and *Gandhi's Truth on the Origins of Militant Nonviolence* (1969), where they are adroitly and sagaciously used by a biographer who is himself a psychiatrist, to Leon Edel's vast biography of Henry James (5 vol., 1953–72), where they are used with sophistication by a man of letters. The science of psychology has also begun to affect the biographer's very approach to his subject: a number of 20th-century authors seek to explore their own involvement with the person they are writing about before embarking upon the life itself.

Ethical. The biographer, particularly the biographer of a contemporary, is often confronted with an ethical problem: how much of the truth, as he has been able to ascertain it, should be printed? Since the inception of biographical criticism in the later 18th century, this somewhat arid—because unanswerable—question has dominated both literary and popular discussion of biographical literature. Upon the publication of the *Life of Samuel Johnson*, James Boswell was bitterly accused of slandering his celebrated subject. More than a century and a half later, Lord Moran's *Winston Churchill: The Struggle for Survival, 1940–1965* (1966), in which Lord Moran used the Boswellian techniques of reproducing conversations from his immediate notes and jottings, was attacked in much the same terms (though the question was complicated by Lord Moran's confidential position as Churchill's physician). In the United States, William Manchester's *Death of a President* (1967), on John F. Kennedy, created an even greater stir in the popular press. There the issue is usually presented as "the public's right to know"; but for the biographer it is a problem of his obligation to preserve historical truth as measured against the personal anguish he may inflict on others in doing so. Since no standard of "biographical morality" has ever been agreed upon—Boswell, Lord Moran, and Manchester have all, for example, had eloquent defenders—the individual biographer must steer his own course. That course in the 20th century is sometimes complicated by the refusal of the custodians of the papers of important persons, particularly national political figures, to provide access to all the documents.

Aesthetic. Biography, while related to history in its search for facts and its responsibility to truth, is truly a branch of literature because it seeks to elicit from facts, by selection and design, the illusion of a life actually being lived. Within the bounds of given data, the biographer seeks to transform plain information into illumination. If he invents or suppresses material in order to create an effect, he fails truth; if he is content to recount facts, he fails art. This tension, between the requirements of authenticity and the necessity for an imaginative ordering of materials to achieve lifelikeness, is perhaps best exemplified in the biographical problem of time. On the one hand, the biographer seeks to portray the unfolding of a life with all its cross-currents of interests, changing emotional states, events; yet in order to avoid reproducing the confusion and clutter of actual daily existence, he must interrupt the flow of diurnal time and group his materials so as to reveal traits of personality, grand themes of experience, and the actions and attitudes leading to moments of high decision. His achievement as a biographical artist will be measured, in great part, by his ability to suggest the sweep of chronology and yet to highlight the major patterns of behaviour that give a life its shape and meaning.

KINDS

Biographies are difficult to classify. It is easily recognizable that there are many kinds of lifewriting, but one kind can easily shade into another; no standard basis for classifi-

Boswell's
great life
of Dr.
Johnson

The
public's
right to
know

Psycholog-
ical inter-
pretation
of
behaviour

cation has yet been developed. A fundamental division offers, however, a useful preliminary view: biographies written from personal knowledge of the subject and those written from research.

Firsthand knowledge. The biography that results from what might be called a vital relationship between the biographer and his subject often represents a conjunction of two main biographical forces: a desire on the part of the writer to preserve "the earthly pilgrimage of a man," as the 19th-century historian Thomas Carlyle calls it (*Critical and Miscellaneous Essays*, 1838), and an awareness that he has the special qualifications, because of direct observation and access to personal papers, to undertake such a task. This kind of biography is, in one form or another, to be found in most of the cultures that preserve any kind of written biographical tradition, and it is commonly to be found in all ages from the earliest literatures to the present. In its first manifestations, it was often produced by, or based upon the recollections of, the disciples of a religious figure—such as the biographical fragments concerning Buddha, portions of the Old Testament, and the Christian gospels. It is sometimes called "source biography" because it preserves original materials, the testimony of the biographer, and often intimate papers of the subject (which have proved invaluable for later biographers and historians—as exemplified by Einhard's 9th-century *Vita Karoli imperatoris* ["Life of Charlemagne"] or Thomas Moore's *Letters and Journals of Lord Byron* [1830]). Biography based on a living relationship has produced a wealth of masterpieces: Tacitus' life of his father-in-law in the *Agricola*, William Roper's life of his father-in-law Sir Thomas More (1626), John Gibson Lockhart's biography (1837–38) of his father-in-law Sir Walter Scott, Johann Peter Eckermann's *Conversations with Goethe* (1836; trans. 1839), and Ernest Jones's *Life and Work of Sigmund Freud* (1953–57). Indeed, what is generally acknowledged as the greatest biography ever written belongs to this class: James Boswell's *Life of Samuel Johnson*.

Research. Biographies that are the result of research rather than firsthand knowledge present a rather bewildering array of forms. First, however, there should be mentioned two special kinds of biographical activity.

Reference collections. Since the late 18th century, the Western world—and, in the 20th century, the rest of the world as well—has produced increasing numbers of compilations of biographical facts concerning both the living and the dead. These collections stand apart from literature. Many nations have multivolume biographical dictionaries such as the *Dictionary of National Biography* in Britain and the *Dictionary of American Biography* in the United States; general encyclopaedias contain extensive information about figures of world importance; classified collections such as *Lives of the Lord Chancellors* (Britain) and biographical manuals devoted to scholars, scientists, and other groups are available in growing numbers; information about living persons is gathered into such national collections as *Who's Who?* (Britain), *Chi è?* (Italy), and *Who's Who in America?*

Character sketches. The short life, however, is a genuine current in the mainstream of biographical literature and is represented in many ages and cultures. Excluding early quasi-biographical materials about religious or political figures, the short biography first appeared in China at about the end of the 2nd century BC, and two centuries later it was a fully developed literary form in the Roman Empire. The *Shih-chi* ("Historical Records"), by Ssu-ma Ch'ien (145?–c. 85 BC), include lively biographical sketches, very short and anecdotal with plentiful dialogue, grouped by character-occupation types such as "maligned statesmen," "rash generals," "assassins," a method that became established tradition with the *Han shu* (*History of the Former Han Dynasty*), by Ssu-ma Ch'ien's successor and imitator, Pan Ku (AD 32–92). Toward the end of the first century AD, in the Mediterranean world, Plutarch's *Lives of the Noble Grecians and Romans*, which are contrasting pairs of biographies, one Greek and one Roman, appeared; there followed within a brief span of years the *Lives of the Caesars*, by the Roman emperor Hadrian's librarian Suetonius. These works established a quite subtle

mingling of character sketch with chronological narrative that has ever since been the dominant mark of this genre. Plutarch, from an ethical standpoint emphasizing the political virtues of man as governor, and Suetonius, from the promptings of sheer biographical curiosity, develop their subjects with telling details of speech and action; and though Plutarch, generally considered to be the superior artist, has greatly influenced other arts than biographical literature—witness Shakespeare's Roman plays, which are based on his *Lives*—Suetonius created in the *Life of Nero* one of the supreme examples of the form. Islāmic literature, from the 10th century, produced short "typed" biographies based on occupation—saints, scholars, and the like—or on arbitrarily chosen personal characteristics. The series of brief biographies has continued to the present day with such representative collections as, in the Renaissance, Giorgio Vasari's *Lives of the Most Eminent Italian Painters, Sculptors, and Architects*, Thomas Fuller's *History of the Worthies of England* in the 17th century, Samuel Johnson's *Lives of the English Poets* in the 18th, and, in more recent times, the "psychographs" of the American Gamaliel Bradford (*Damaged Souls*, 1923), Lytton Strachey's *Eminent Victorians* (1918) and the "profiles" that have become a hallmark of the weekly magazine *The New Yorker*.

Further classification of biographies compiled by research can be achieved by regarding the comparative objectivity of approach. For convenience, six categories, blending one into the other in infinite gradations and stretching from the most objective to the most subjective, can be employed.

Informative biography. This, the first category, is the most objective and is sometimes called "accumulative" biography. The author of such a work, avoiding all forms of interpretation except selection—for selection, even in the most comprehensive accumulation, is inevitable—seeks to unfold a life by presenting, usually in chronological order, the paper remains, the evidences, relating to that life. This biographer takes no risks but, in turn, seldom wins much critical acclaim: his work is likely to become a prime source for biographers who follow him. During the 19th century, the *Life of Milton: Narrated in Connection with the Political, Ecclesiastical, and Literary History of his Time* (7 vol., 1859–94), by David Masson, and *Abraham Lincoln: A History* (10 vol., 1890), by John G. Nicolay and John Hay, offer representative samples. In the 20th century such works as Edward Nehls's, *D.H. Lawrence: A Composite Biography* (1957–59) and David Alec Wilson's collection of the life records of Thomas Carlyle (1923–29), in six volumes, continue the traditions of this kind of life writing.

Critical biography. This second category, scholarly and critical, unlike the first, does offer a genuine presentation of a life. These works are very carefully researched; sources and "justifications" (as the French call them) are scrupulously set forth in notes, appendixes, bibliographies; inference and conjecture, when used, are duly labelled as such; no fictional devices or manipulations of material are permitted, and the life is generally developed in straight chronological order. Yet such biography, though not taking great risks, does employ the arts of selection and arrangement. The densest of these works, completely dominated by fact, have small appeal except to the specialist. Those written with the greatest skill and insight are in the first rank of modern life writing. In these scholarly biographies—the "life and times" or the minutely detailed life—the author is able to deploy an enormous weight of matter and yet convey the sense of a personality in action, as exemplified in Leslie Marchand's *Byron* (1957), with some 1,200 pages of text and 300 pages of notes, Dumas Malone's *Jefferson and his Time* (4 vol., 1948–70), Churchill's *Marlborough* (1933–38), Douglas S. Freeman's *George Washington* (1948–57). The critical biography aims at evaluating the works as well as unfolding the life of its subject, either by interweaving the life in its consideration of the works or else by devoting separate chapters to the works. Critical biography has had its share of failures: except in skillful hands, criticism clumsily intrudes upon the continuity of a life, or the works of the subject are made to yield doubtful interpretations of character, particularly

The relationship between writer and subject

Inevitability of selection in biography

The first short biography

The problems of critical biography

in the case of literary figures. It has to its credit, however, such fine biographies as Arthur S. Link, *Wilson* (5 vol., 1947–65); Richard Ellmann, *James Joyce* (1959); Ernest Jones, *The Life and Works of Sigmund Freud*; Douglas S. Freeman, *Lee* (1934–35); and Edgar Johnson, *Charles Dickens* (1952).

"Standard" biography. This third, and central, category of biography, balanced between the objective and the subjective, represents the mainstream of biographical literature, the practice of biography as an art. From antiquity until the present—within the limits of the psychological awareness of the particular age and the availability of materials—this kind of biographical literature has had as its objective what Sir Edmund Gosse called "the faithful portrait of a soul in its adventures through life." It seeks to transform, by literary methods that do not distort or falsify, the truthful record of fact into the truthful effect of a life being lived. Such biography ranges in style and method from George Cavendish's 16th-century life of Cardinal Wolsey, Roger North's late-17th-century lives of his three brothers, and Boswell's life of Johnson to modern works like Lord David Cecil's *Melbourne*, Garrett Mattingly's *Catherine of Aragon*, Andrew Turnbull's *Scott Fitzgerald*, and Leon Edel's *Henry James*.

Interpretative biography. This fourth category of life writing is subjective and has no standard identity. At its best it is represented by the earlier works of Catherine Drinker Bowen, particularly her lives of Tchaikovsky, "*Beloved Friend*" (1937), and Oliver Wendell Holmes, *Yankee from Olympus* (1944). She molds her sources into a vivid narrative, worked up into dramatic scenes that always have some warranty of documentation—the dialogue, for example, is sometimes devised from the indirect discourse of letter or diary. She does not invent materials; but she quite freely manipulates them—that is to say, interprets them—according to the promptings of insight, derived from arduous research, and with the aim of unfolding her subject's life as vividly as possible. (Mrs. Bowen, much more conservative in her later works, clearly demonstrates the essential distance between the third and fourth categories: her distinguished life of Sir Edward Coke, *The Lion and the Throne* [1957], foregoes manipulation and the "recreation" of dialogue and limits interpretation to the artful deployment of biographical resources.) Very many interpretative biographies stop just short of fictionalizing in the freedom with which they exploit materials. The works of Frank Harris (*Oscar Wilde*, 1916) and Hesketh Pearson (*Tom Paine, Friend of Mankind*, 1937; *Beerbohm Tree*, 1956) demonstrate this kind of biographical latitude.

Fictionalized biography. The books in this fifth category belong to biographical literature only by courtesy. Materials are freely invented, scenes and conversations are imagined; unlike the previous category, this class often depends almost entirely upon secondary sources and cursory research. Its authors, well represented on the paperback shelves, have created a hybrid form designed to mate the appeal of the novel with a vague claim to authenticity. This form is exemplified by writers such as Irving Stone, in his *Lust for Life* (on van Gogh) and *The Agony and the Ecstasy* (on Michelangelo). Whereas the compiler of biographical information (the first category) risks no involvement, the fictionalizer admits no limit to it.

Fiction presented as biography. The sixth and final category is outright fiction, the novel written as biography or autobiography. It has enjoyed brilliant successes. Such works do not masquerade as lives; rather, they imaginatively take the place of biography where perhaps there can be no genuine life writing for lack of materials. Among the most highly regarded examples of this genre are, in the guise of autobiography, Robert Graves's books on the Roman emperor Claudius, *I, Claudius* and *Claudius the God and His Wife Messalina*; Mary Renault's *The King Must Die* on the legendary hero Theseus; and Marguerite Yourcenar's *Memoirs of Hadrian*. The diary form of autobiography was amusingly used by George and Weedon Grossmith to tell the trials and tribulations of their fictional character, Charles Poster, in *The Diary of a Nobody* (1892). In the form of biography this category includes Graves's *Count Belisarius* and Hope Muntz's *Golden War-*

rior (on Harold II, vanquished at the Battle of Hastings, 1066). Some novels-as-biography, using fictional names, are designed to evoke rather than re-create an actual life, such as W. Somerset Maugham's *Moon and Sixpence* (Gauguin) and *Cakes and Ale* (Thomas Hardy) and Robert Penn Warren's *All the King's Men* (Huey Long).

"Special-purpose" biography. In addition to these six main categories, there exists a large class of works that might be denominated "special-purpose" biography. In these works the art of biography has become the servant of other interests. They include potboilers (written as propaganda or as a scandalous exposé) and "as-told-to" narratives (often popular in newspapers) designed to publicize a celebrity. This category includes also "campaign biographies" aimed at forwarding the cause of a political candidate (Nathaniel Hawthorne's *Life of Franklin Pierce* [1852] being an early example); the weighty commemorative volume, not infrequently commissioned by the widow (which, particularly in Victorian times, has usually enshrouded the subject in monotonous eulogy); and pious works that are properly called hagiography, or lives of holy men, written to edify the reader.

Informal autobiography. Autobiography, like biography, manifests a wide variety of forms, beginning with the intimate writings made during a life that were not intended (or apparently not intended) for publication.

Letters, diaries, and journals. Broadly speaking, the order of this category represents a scale of increasingly self-conscious revelation. Collected letters, especially in carefully edited modern editions such as W.S. Lewis' of the correspondences of the 18th-century man of letters Horace Walpole (34 vol., 1937–65), can offer a rewarding though not always predictable experience: some eminent people commit little of themselves to paper, while other lesser figures pungently re-create themselves and their world. The 15th-century *Paston Letters* constitute an invaluable chronicle of the web of daily life woven by a tough and vigorous English family among the East Anglian gentry during the Wars of the Roses; the composer Mozart and the poet Byron, in quite different ways, are among the most revealing of letter writers. Diarists have made great names for themselves out of what seems a humble branch of literature. To mention only two, in the 20th century the young Jewish girl Anne Frank created such an impact by her recording of narrow but intense experience that her words were translated to stage and screen; while a comparatively minor figure of 17th-century England, Samuel Pepys—he was secretary to the navy—has immortalized himself in a diary that exemplifies the chief qualifications for this kind of writing—candour, zest, and an unself-conscious enjoyment of self. The somewhat more formal journal is likewise represented by a variety of masterpieces, from the notebooks, which reveal the teeming, ardent brain of Leonardo da Vinci, and William Wordsworth's sister Dorothy's sensitive recording of experience in her *Journals* (1897), to French foreign minister Armand de Caulaincourt's recounting of his flight from Russia with Napoleon (translated as *With Napoleon in Russia*, 1935) and the *Journals* of the brothers Goncourt, which present a confidential history of the literary life of mid-19th-century Paris.

Memoirs and reminiscences. These are autobiographies that usually emphasize *what* is remembered rather than *who* is remembering; the author, instead of recounting his life, deals with those experiences of his life, people, and events that he considers most significant. (The extreme contrast to memoirs is the spiritual autobiography, so concentrated on the life of the soul that the author's outward life and its events remains a blur. The artless *res gestae*, a chronology of events, occupies the middle ground.)

In the 15th century, Philippe de Commines, modestly effacing himself except to authenticate a scene by his presence, presents in his *Mémoires* a life of Louis XI, master of statecraft, as witnessed by one of the most sagacious counsellors of the age. The memoirs of Giacomo Casanova boast of an 18th-century rake's adventures; those of Hector Berlioz explore with great brilliance the trials of a great composer, the reaches of an extraordinary personality, and the musical life of Europe in the first part

Cursory
research
in fiction-
alized
biography

The nature
of memoirs

of the 19th century. The memoir form is eminently represented in modern times by Sir Osbert Sitwell's polished volumes, presenting a tapestry of recollections that, as has been observed, "tells us little about what it feels like to be in Sir Osbert's skin"—a phrase perfectly illustrating the difference between memoirs and formal autobiography.

Formal autobiography. This category offers a special kind of biographical truth: a life, reshaped by recollection, with all of recollection's conscious and unconscious omissions and distortions. The novelist Graham Greene says that, for this reason, an autobiography is only "a sort of life" and uses the phrase as the title for his own autobiography (1971). Any such work is a true picture of what, at one moment in a life, the subject wished—or is impelled—to reveal of that life. An event recorded in the autobiographer's youthful journal is likely to be somewhat different from that same event recollected in later years. Memory being plastic, the autobiographer regenerates his materials as he uses them. The advantage of possessing unique and private information, accessible to no researching biographer, is counterbalanced by the difficulty of establishing a stance that is neither overmodest nor aggressively self-assertive. The historian Edward Gibbon declares, "... I must be conscious that no one is so well qualified as myself to describe the service of my thoughts and actions." The 17th-century English poet Abraham Cowley provides a rejoinder: "It is a hard and nice subject for a man to write of himself; it grates his own heart to say anything of disparagement and the reader's ears to hear anything of praise from him."

There are but few and scattered examples of autobiographical literature in antiquity and the Middle Ages. In the 2nd century BC the Chinese classical historian Ssu-ma Ch'ien included a brief account of himself in the *Shihchi*, "Historical Records." It is stretching a point to include, from the 1st century BC, the letters of Cicero (or, in the early Christian era, the letters of St. Paul); and Julius Caesar's *Commentaries* tell little about Caesar, though they present a masterly picture of the conquest of Gaul and the operations of the Roman military machine at its most efficient. The *Confessions* of St. Augustine, of the 5th century AD, belong to a special category of autobiography discussed below; the 14th-century *Letter to Posterity* of the Italian poet Petrarch is but a brief excursion in the field.

Speaking generally, then, it can be said that autobiography begins with the Renaissance in the 15th century; and, surprisingly enough, the first example was written not in Italy but in England by a woman entirely untouched by the "new learning" or literature. In her old age Margery Kempe, the sobbing mystic, or hysteric, of Lynn in Norfolk, dictated an account of her bustling, far-faring life, which, however concerned with religious experience, racily reveals her somewhat abrasive personality and the impact she made upon her fellows. This is done in a series of scenes, mainly developed by dialogue. Though calling herself, in abject humility, "the creature," Margery knew, and has effectively transmitted the proof, that she was a remarkable person.

The first full-scale formal autobiography was written a generation later by a celebrated Humanist publicist of the age, Enea Silvio Piccolomini, after he was elevated to the papacy, in 1458, as Pius II—the result of an election that he recounts with astonishing frankness spiced with malice. In the first book of his autobiography—misleadingly named *Commentarii*, in evident imitation of Caesar—Pius II traces his career up to becoming pope; the succeeding 11 books (and a fragment of a 12th, which breaks off a few months before his death in 1464) present a panorama of the age, with its cruel and cultivated Italian tyrants, cynical *condottieri* (professional soldiers), recalcitrant kings, the politics and personalities behind the doors of the Vatican, and the urbane but exuberant character of the Pope himself. Pius II exploits the plasticity of biographical art by creating opportunities—especially when writing of himself as the connoisseur of natural beauties and antiquities—for effective autobiographical narration. His "Commentaries" show the art of formal autobiography in full bloom in its beginnings; they rank as one of its half dozen greatest exemplars.

The neglected autobiography of the Italian physician and astrologer Gironimo Cardano, a work of great charm, and the celebrated adventures of the goldsmith and sculptor Benvenuto Cellini in Italy of the 16th century; the uninhibited autobiography of the English historian and diplomat Lord Herbert of Cherbury, in the early 17th; and Colley Cibber's *Apology for The Life of Colley Cibber, Comedian* in the early 18th—these are representative examples of biographical literature from the Renaissance to the Age of Enlightenment. The latter period itself produced three works that are especially notable for their very different reflections of the spirit of the times as well as of the personalities of their authors: the urbane autobiography of Edward Gibbon, the great historian; the plainspoken, vigorous success story of an American who possessed all the talents, Benjamin Franklin; and the somewhat morbid introspection of a revolutionary Swiss-French political and social theorist, the *Confessions* of J.-J. Rousseau—the latter leading to two autobiographical explorations in poetry during the Romantic Movement (flourished 1798–1837) in England, Wordsworth's *Prelude* and Byron's *Childe Harold*, cantos III and IV. Significantly, it is at the end of the 18th century that the word autobiography apparently first appears in print, in *The Monthly Review*, 1797.

Specialized forms. These might roughly be grouped under four heads: thematic, religious, intellectual, and fictionalized. The first grouping includes books with such diverse purposes as Adolf Hitler's *Mein Kampf* (1924), *The Americanization of Edward Bok* (1920), and Richard Wright's *Native Son* (1940). Religious autobiography claims a number of great works, ranging from the *Confessions* of St. Augustine and Peter Abelard's *Historia Calamitatum* (*The Story of My Misfortunes*) in the Middle Ages to the autobiographical chapters of Thomas Carlyle's *Sartor Resartus* ("The Everlasting No," "Centre of Indifference," "The Everlasting Yea") and Cardinal John Newman's beautifully wrought *Apologia* in the 19th century. That century and the early 20th saw the creation of several intellectual autobiographies. The *Autobiography* of the philosopher John S. Mill, severely analytical, concentrates upon "an education which was unusual and remarkable." It is paralleled, across the Atlantic, in the bleak but astringent quest of *The Education of Henry Adams* (printed privately 1906; published 1918). Edmund Gosse's sensitive study of the difficult relationship between himself and his Victorian father, *Father and Son* (1907), and George Moore's quasi-novelized crusade in favour of Irish art, *Hail and Farewell* (1911–14), illustrate the variations of intellectual autobiography. Finally, somewhat analogous to the novel as biography (for example, Graves's *I, Claudius*) is the autobiography thinly disguised as, or transformed into, the novel. This group includes such works as Samuel Butler's *Way of All Flesh* (1903), James Joyce's *Portrait of the Artist as a Young Man* (1916), George Santayana's *Last Puritan* (1935), and the gargantuan novels of Thomas Wolfe (*Look Homeward, Angel* [1929], *Of Time and the River* [1935]).

HISTORICAL DEVELOPMENT

Western literature. *Antiquity.* In the Western world, biographical literature can be said to begin in the 5th century BC with the poet Ion of Chios, who wrote brief sketches of such famous contemporaries as Pericles and Sophocles. It continued throughout the classical period for a thousand years, until the dissolution of the Roman Empire in the 5th century AD. Broadly speaking, the first half of this period exhibits a considerable amount of biographical activity, of which much has been lost; such fragments as remain of the rest—largely funeral elegies and rhetorical exercises depicting ideal types of character or behaviour—suggest that from a literary point of view the loss is not grievous. (An exception is the life of the Roman art patron Pomponius Atticus, written in the 1st century BC by Cornelius Nepos.) Biographical works of the last centuries in the classical period, characterized by numerous sycophantic accounts of emperors, share the declining energies of the other literary arts. But although there are few genuine examples of life writing, in the modern sense of the term, those few are masterpieces. The two greatest teachers

Plato's
accounts of
Socrates

of the classical Mediterranean world, Socrates and Jesus Christ, both prompted the creation of magnificent biographies written by their followers. To what extent Plato's life of Socrates keeps to strict biographical truth cannot now be ascertained (though the account of Socrates given by Plato's contemporary the soldier Xenophon, in his *Memorabilia*, suggests a reasonable faithfulness) and he does not offer a full-scale biography. Yet in his two consummate biographical dialogues—*The Apology* (recounting the trial and condemnation of Socrates) and the *Phaedo* (a portrayal of Socrates' last hours and death)—he brilliantly recreates the response of an extraordinary character to the crisis of existence. Some 400 years later there came into being four lives of Jesus, the profound religious significance of which has inevitably obscured their originality—their homely detail, anecdotes, and dialogue that, though didactic in purpose, also evoke a time and a personality. The same century, the first of the Christian era, gave birth to the three first truly "professional" biographers—Plutarch and Suetonius (discussed above) and the historian Tacitus, whose finely wrought biography of his father-in-law, *Agricola*, concentrating on the administration rather than the man, has something of the monumental quality of Roman architecture. The revolution in thought and attitude brought about by the growth of Christianity is signalled in a specialized autobiography, the *Confessions* of St. Augustine; but the biographical opportunity suggested by Christian emphasis on the individual soul was, oddly, not to be realized. If the blood of the martyrs fertilized the seed of the new faith, it did not promote the art of biography. The demands of the church and the spiritual needs of men, in a twilight world of superstition and violence, transformed biography into hagiography. There followed a thousand years of saints' lives: the art of biography forced to serve ends other than its own.

Middle Ages. This was a period of biographical darkness, an age dominated by the priest and the knight. The priest shaped biography into an exemplum of otherworldliness, while the knight found escape from daily brutishness in allegory, chivalric romances, and broad satire (the *fabliaux*). Nevertheless, glimmerings can be seen. A few of the saints' lives, like Eadmer's *Life of Anselm*, contain anecdotal materials that give some human flavour to their subjects; the 13th-century French nobleman Jean, sire de Joinville's life of St. Louis (Louis IX of France), *Mémoires*, offers some lively scenes. The three most interesting biographical manifestations came early. Bishop Gregory of Tours' *History of the Franks* depicts artlessly but vividly, from firsthand observation, the lives and personalities of the four grandsons of Clovis and their fierce queens in Merovingian Gaul of the 6th century. Bede's *Ecclesiastical History of the English People*, of the 8th century, though lacking the immediacy and exuberance—and the violent protagonists—of Gregory, presents some valuable portraits, like those of "the little dark man," Paulinus, who converted the King of Northumbria to Christianity.

Einhard's
Life of Charlemagne

Most remarkable, however, a self-consciously wrought work of biography came into being in the 9th century: this was *The Life of Charlemagne*, written by a cleric at his court named Einhard. He is aware of his biographical obligations and sets forth his point of view and his motives:

I have been careful not to omit any facts that could come to my knowledge, but at the same time not to offend by a prolix style those minds that despise everything modern . . . No man can write with more accuracy than I of events that took place about me, and of facts concerning which I had personal knowledge. . . .

He composes the work in order to ensure that Charlemagne's life is not "wrapped in the darkness of oblivion" and out of gratitude for "the care that King Charles bestowed upon me in my childhood, and my constant friendship with himself and his children." Though Einhard's biography, by modern standards, lacks sustained development, it skillfully reveals the chief patterns of Charlemagne's character—his constancy of aims, powers of persuasion, passion for education. Einhard's work is far closer to modern biography than the rudimentary poetry and drama of his age are to their modern counterparts.

Renaissance. Like the other arts, biography stirs into fresh life with the Renaissance in the 15th century. Its most significant examples were autobiographical, as has already been mentioned. Biography was chiefly limited to uninspired panegyrics of Italian princes by their court Humanists, such as Simonetta's life of the great *condottiere*, Francesco Sforza, duke of Milan.

During the first part of the 16th century in England, now stimulated by the "new learning" of Erasmus, John Colet, Thomas More, and others, there were written three works that can be regarded as the initiators of modern biography: More's *History of Richard III*, William Roper's *Mirroure of Vertue in Worldly Greatness; or, the life of Syr Thomas More*, and George Cavendish's *Life of Cardinal Wolsey*. The *History of Richard III* (written about 1513 in both an English and a Latin version) unfortunately remains unfinished; and it cannot meet the strict standards of biographical truth since, under the influence of classical historians, a third of the book consists of dialogue that is not recorded from life. However, it is a brilliant work, exuberant of wit and irony, that not only constitutes a biographical landmark but is also the first piece of modern English prose. With relish, More thus sketches Richard's character:

He was close and secret, a deep dissembler, lowly of countenance, arrogant of heart, outwardly companionable where he inwardly hated, not hesitating to kiss whom he thought to kill.

Worked up into dramatic scenes, this biography, as reproduced in the *Chronicles* of Edward Hall and Raphael Holinshed, later provided both source and inspiration for Shakespeare's rousing melodramatic tragedy, *Richard III*. The lives written by Roper and Cavendish display interesting links, though the two men were not acquainted: they deal with successive first ministers destroyed by that brutal master of politics, Henry VIII; they are written from first hand observation of their subjects by, respectively, a son-in-law and a household officer; and they exemplify, though never preach, a typically Renaissance theme: *Indignatio principis mors est*—"the Prince's anger is death." Roper's work is shorter, more intimate, and simpler; in a series of moving moments it unfolds the struggle within Sir Thomas More between his duty to conscience and his duty to his king. Cavendish offers a more artful and richly developed narrative, beautifully balanced between splendid scenes of Wolsey's glory and vanity and ironically contrasting scenes of disgrace, abasement, and painfully achieved self-knowledge.

Biography
in the
Renaissance

The remaining period of the Renaissance, however, is disappointingly barren. In Russia, where medieval saints' lives had also been produced, there appears a modest biographical manifestation in the *Stepennaya Kniga* ("Book of Degrees," 1563), a collection of brief lives of princes and prelates. Somewhat similarly, in France, which was torn by religious strife, Pierre Brantôme wrote his *Lives of Famous Ladies* and *Lives of Famous Men*. The Elizabethan Age in England, for all its magnificent flowering of the drama, poetry, and prose, did not give birth to a single biography worthy of the name. Sir Fulke Greville's account of Sir Philip Sidney (1652) is marred by tedious moralizing; Francis Bacon's accomplished life of the first Tudor monarch, *The Historie of the Raigne of King Henry the Seventh* (1622), turns out to be mainly a history of the reign. But Sir Walter Raleigh suggests an explanation for this lack of biographical expression in the introduction to his *History of the World* (1614): "Whosoever, in writing a modern history, shall follow truth too near the heels, it may haply strike out his teeth"—as Sir John Hayward could testify, having been imprisoned in the Tower of London because his account (1599) of Richard II's deposition, two centuries earlier, had aroused Queen Elizabeth's anger.

17th and 18th centuries. In the 17th century the word biography was first employed to create a separate identity for this type of writing. That century and the first half of the 18th presents a busy and sometimes bizarre biographical landscape. It was an era of experimentation and preparation rather than of successful achievement. In the New World, the American Colonies began to develop a scattered biographical activity, none of it of lasting impor-

tance, France offers the celebrated *Letters* of the Marquise de Sévigné to her daughter, an intimate history of the Age of Louis XIV; numerous memoirs, such as those of Louis de Rouvroy, duc de Saint-Simon, and the acerbic ones of the Cardinal de Retz (1717); and the philosopher and critic Pierre Bayle's *Dictionnaire historique et critique* (1697), which was followed by specialized biographical collections and reference works. England saw an outpouring, beginning in the earlier 17th century, of Theophrastan "characters" (imaginary types imitated from the work of Theophrastus, a follower of Aristotle), journals, diaries, the disorganized but vivid jottings of John Aubrey (later published in 1898 as *Brief Lives*); and in the earlier 18th century there were printed all manner of sensational exposés, biographical sketches of famous criminals, and the like. In this era women appear for the first time as biographers. Lady Fanshawe wrote a life of her ambassador-husband (1829); Lucy Hutchinson, one of her Puritan warrior-husbands (written after 1664, published 1806); and Margaret Cavendish, duchess of Newcastle, produced a warm, bustling life—still good reading today—of her duke, an amiable mediocrity (*The Life of the thrice Noble Prince William Cavendish, Duke Marquess, and Earl of Newcastle*, 1667). This age likewise witnessed the first approach to a professional biographer, the noted lover of angling, Izaak Walton, whose five lives (of the poets John Donne [1640] and George Herbert [1670], the diplomat Sir Henry Wotton [1651], and the ecclesiastics Richard Hooker [1665] and Robert Sanderson [1678]) tend to endow their diverse subjects with something of Walton's own genteel whimsicality but nonetheless create skillful biographical portraits. The masterpieces of the age are unquestionably Roger North's biographies (not published until 1742, 1744) of his three brothers: Francis, the lord chief justice, "my best brother"; the lively merchant-adventurer Sir Dudley, his favourite; and the neurotic scholar John. Also the author of an autobiography, Roger North likewise produced, as a preface to his life of Francis, the first extensive critical essay on biography, which anticipates some of the ideas of Samuel Johnson and James Boswell.

The last half of the 18th century witnessed the remarkable conjunction of these two remarkable men, from which sprang what is generally agreed to be the world's supreme biography, Boswell's *Life of Samuel Johnson LL.D.* (1791). Dr. Johnson, literary dictator of his age, critic and lexicographer who turned his hand to many kinds of literature, himself created the first English professional biographies in *The Lives of the English Poets*. In essays and in conversation, Johnson set forth principles for biographical composition: the writer must tell the truth—"the business of the biographer is often to . . . display the minute details of daily life," for it is these details that recreate a living character; and men need not be of exalted fame to provide worthy subjects.

For more than one reason the somewhat disreputable and incredibly diligent Scots lawyer James Boswell can be called the unique genius of biographical literature, bestriding both autobiography and biography. Early in his acquaintance with Johnson he was advised by the Doctor "to keep a journal of my [Boswell's] life, full and unreserved." Boswell followed this advice to the letter. His gigantic journals offer an unrivalled self-revelation of a fascinatingly checkered character and career—whether as a young rake in London or thrusting himself upon the aged Rousseau or making his way to Voltaire's seclusion at Ferney in Switzerland with the aim of converting that celebrated skeptic to Christianity. Boswell actively helped to stage the life of Johnson that he knew he was going to write—drawing out Johnson in conversation, setting up scenes he thought likely to yield rich returns—and thus, at moments, he achieved something like the novelist's power over his materials, being himself an active part of what he was to re-create. Finally, though he invented no new biographical techniques, in his *Life of Samuel Johnson* he interwove with consummate skill Johnson's letters and personal papers, Johnson's conversation as assiduously recorded by the biographer, material drawn from interviews with large numbers of people who knew Johnson, and his own observation of Johnson's behaviour, to elicit

the living texture of a life and a personality. Boswell makes good his promise that Johnson "will be seen as he really was . . ." The influence of Boswell's work penetrated throughout the world and, despite the development of new attitudes in biographical literature, has persisted to this day as a pervasive force. Perhaps equally important to life writers has been the inspiration provided by the recognition accorded Boswell's *Life* as a major work of literary art. Since World War II there have often been years, in the United States, when the annual bibliographies reveal that more books or articles were published about Johnson and Boswell than about all the rest of biographical literature together.

19th century. The *Life of Johnson* may be regarded as a representative psychological expression of the Age of Enlightenment, and it certainly epitomizes several typical characteristics of that age: devotion to urban life, confidence in common sense, emphasis on man as a social being. Yet in its extravagant pursuit of the life of one individual, in its laying bare the eccentricities and suggesting the inner turmoil of personality, it may be thought of as part of that revolution in self-awareness, ideas, aspirations, exemplified in Rousseau's *Confessions*, the French Revolution, the philosophical writings of the German philosopher Immanuel Kant, the political tracts of Thomas Paine, and the works of such early Romantic poets as Robert Burns, William Blake, Wordsworth—a revolution that in its concern with the individual psyche and the freedom of man seemed to augur well for biographical literature. This promise, however, was not fulfilled in the 19th century.

That new nation, the United States of America, despite the stimulus of a robust and optimistic society, flamboyant personalities on the frontier, a generous share of genius, and the writing of lives by eminent authors such as Washington Irving and Henry James, produced no biographies of real importance. One professional biographer, James Parton, published competent, well-researched narratives, such as his lives of Aaron Burr and Andrew Jackson, but they brought him thin rewards and are today outmoded. In France, biography was turned inward, to romantic introspection, a trend introduced by Étienne Pivert de Senancour's *Obermann* (1804). It was followed by autobiographies thinly disguised as novels such as Benjamin Constant's *Adolphe* (1816), *La Vie de Henri Brulard* of Stendhal (Marie Henri Beyle), and similar works by Alphonse de Lamartine and Alfred de Musset, in which the emotional malaise of the hero is subjected to painstaking analysis. In Great Britain the 19th century opened promisingly with an outburst of biographical-autobiographical production, much of which came from prominent figures of the Romantic Movement, including Samuel Taylor Coleridge, Robert Southey, William Hazlitt, and Thomas De Quincey. Thomas Moore's *Letters and Journals of Lord Byron* (1830), John Gibson Lockhart's elaborate life (1837–38) of his father-in-law, Sir Walter Scott, and, later, Elizabeth Cleghorn Gaskell's *Life of Charlotte Brontë* (1857), James Anthony Froude's study of Carlyle (2 vol. 1882; 2 vol. 1884), John Forster's *Life of Charles Dickens* (1872–74) all followed, to some degree, what may loosely be called the Boswell formula. Yet most of these major works are marred by evasions and omissions of truth—though Lockhart and Froude, for example, were attacked as conscienceless despoilers of the dead—and, before the middle of the century, biography was becoming stifled. As the 20th century biographer and critic Sir Harold Nicolson wrote in *The Development of English Biography* (1927), "Then came earnestness, and with earnestness hagiography descended on us with its sullen cloud . . ." Insistence on respectability, at the expense of candour, had led Carlyle to observe acridly, "How delicate, how decent is English biography, bless its mealy mouth!" and to pillory its productions as "vacuum-biographies."

20th century. The period of modern biography was ushered in, generally speaking, by World War I. All the arts were in ferment, and biographical literature shared in the movement, partly as a reaction against 19th-century conventions, partly as a response to advances in psychology, and partly as a search for new means of expression. This

The first
"professional"
biographer

Boswell's
influence
on
biography

Biography
during the
Victorian
age

revolution, unlike that at the end of the 18th century, was eventually destined to enlarge and enhance the stature of biography. The chief developments of modern life writing may be conveniently classified under five heads: (1) an increase in the numbers and general competence of biographies throughout the Western world; (2) the influence on biographical literature of the counterforces of science and fictional writing; (3) the decline of formal autobiography and of biographies springing from a personal relationship; (4) the range and variety of biographical expression; and (5) the steady, though moderate, growth of a literature of biographical criticism. Only the first three of these developments need much elaboration.

Little has been said about biography since the Renaissance in Germany, Spain, Italy, Scandinavia, and the Slavic countries because, as in the case of Russia, there had been comparatively little biographical literature and because biographical trends, particularly since the end of the 18th century, generally followed those of Britain and France. Russian literary genius in prose is best exemplified during both the 19th and 20th centuries in the novel. In the 19th century, however, Leo Tolstoy's numerous autobiographical writings, such as *Childhood* and *Boyhood*, and Sergey Aksakov's *Years of Childhood* and *A Russian Schoolboy*, and in the 20th century, Maksim Gorky's autobiographical trilogy (*Childhood: In the World*; and *My Universities*, 1913–23) represent, in specialized form, a limited biographical activity. The close control of literature exercised by the 20th-century Communist governments of eastern Europe has created a wintry climate for biography. The rest of Europe, outside the iron curtain, has manifested in varying degrees the fresh biographical energies and practices illustrated in British-American life writing: biography is now, as never before, an international art that shares a more or less common viewpoint.

The second characteristic of modern biography, its being subject to the opposing pressures of science and fictional writing, has a dark as well as a bright side. Twentieth-century fiction, boldly and restlessly experimental, has, on the one hand, influenced the biographer to aim at literary excellence, to employ devices of fiction suitable for biographical ends; but, on the other, fiction has also probably encouraged the production of popular pseudo-biography, hybrids of fact and fancy, as well as of more subtle distortions of the art form. Science has exerted two quite different kinds of pressure: the prestige of the traditional sciences, in their emphasis on exactitude and rigorous method, has undoubtedly contributed to a greater diligence in biographical research and an uncompromising scrutiny of evidences; but science's vast accumulating of facts—sometimes breeding the worship of fact for its own sake—has helped to create an atmosphere in which today's massive, note-ridden and fact-encumbered lives proliferate and has probably contributed indirectly to a reluctance in the scholarly community to take the risks inevitable in true biographical composition.

The particular science of psychology, as earlier pointed out, has conferred great benefits upon the responsible practitioners of biography. It has also accounted in large part, it would appear, for the third characteristic of modern biography: the decline of formal autobiography and of the grand tradition of biography resulting from a personal relationship. For psychology has rendered the self more exposed but also more elusive, more fascinatingly complex and, in the darker reaches, somewhat unpalatable. Since honesty would force the autobiographer into a self-examination both formidable to undertake and uncomfortable to publish, instead he generally turns his attention to outward experiences and writes memoirs and reminiscences—though France offers something of an exception in the journals of such writers as André Gide (1947–51), Paul Valéry (1957), François Mauriac (1934–50), Julien Green (1938–58). Similarly, psychology, in revealing the fallacies of memory, the distorting power of an emotional relationship, the deceits of observation, has probably discouraged biography written by a friend of its subject. Moreover, so many personal papers are today preserved that a life-long friend of the subject scarcely has time to complete his biography.

After World War I, the work of Lytton Strachey played a somewhat similar role to that of Boswell in heading a "revolution" in biography. *Eminent Victorians* and *Queen Victoria* (1921), followed by *Elizabeth and Essex* (1928), with their artful selection, lacquered style, and pervasive irony, exerted an almost intoxicating influence in the 1920s and '30s. Writers seeking to capitalize on Strachey's popularity and ape Strachey's manner, without possessing Strachey's talents, produced a spate of "debunking" biographies zestfully exposing the clay feet of famous historical figures. By World War II, however, this kind of biography had been discredited; Strachey's adroit detachment and literary skill were recognized to be his true value, not his dangerously interpretative method; and, since that time, biography has steadied into an established, if highly varied, form of literature.

Other literatures. Biography as an independent art form, with its concentration upon the individual life and its curiosity about the individual personality, is essentially a creation of Western man. In the Orient, for all its long literary heritage, and in Islām, too, biographical literature does not show the development, nor assume the importance, of Western life writing. In China, until comparatively recently, biography had been an appendage, or by-product, of historical writing and scholarly preoccupation with the art of government, in the continuing tradition of the "Historical Records" of Ssu-ma Ch'ien and Pan Ku. In India it has been the enduring concern for spiritual values and for contemplation or mystical modes of existence that have exerted the deepest influence on literature from the first millennium BC to the present, and this has not provided a milieu suitable to biographical composition. Generally speaking, the literary history of Japan, too, offers only fragmentary or limited examples of life writing.

It was not until the beginning of the 20th century in China that biography began to appear as an independent form (and this was evidently the result of western influence), when Liang Ch'i-ch'ao (1873–1929) wrote a number of lives, including one of Confucius, and was followed by Hu Shih (1891–1962), who, like his predecessor, worked to promote biographical composition as an art form. Except for China after the establishment of the Communist state in 1949, biography in the Orient—notably in India and Japan—has shared, to a limited extent, the developments in biographical literature demonstrated in the rest of the world.

BIOGRAPHICAL LITERATURE TODAY

In the United States, Great Britain, and the rest of the Western world generally, biography today enjoys a moderate popular and critical esteem. In the year 1929, at the height of the biographical "boom," there were published in the United States 667 new biographies; in 1962 exactly the same number appeared, the population in the meantime having increased by something like 50 percent. On the average, in the English-speaking world, biographical titles account for approximately 5 percent of the annual output of books. Yet they have won their share of literary prizes and for their authors a considerable degree of literary eminence; if few universally acclaimed masterpieces are being produced, it is probably true that the art of biography is seeing a higher general level of achievement than ever before. The recreation of a life is also now being attempted in other media than that of prose. Biographical drama has of course been staged from before the time of Shakespeare; it continues to be popular, whether translated from narrative to the theatre (as the *Diary of Anne Frank*) or written specifically for the stage, like Jean Anouilh's *Becket* and Robert Bolt's study of Sir Thomas More, *A Man for All Seasons* (which nonetheless owes a great deal to William Roper). The cinema often follows with its versions of such plays; it likewise produces original biographical films, generally with indifferent success. Television, too, offers historical "recreations" of various sorts, and with varying degrees of responsibility, but has achieved only a few notable examples of biographical illumination, for the conflict between gripping visual presentation and the often undramatic, but important, biographical truth is difficult to resolve. Biography, indeed, seems less innovative, less

rewarding of experiment, and less adaptable to new media, than does fiction or perhaps even history. Words are no longer the only way to tell a story and perhaps in time will not be regarded as the chief way; but so far they seem the best way of unfolding the full course of a life and exploring the quirks and crannies of a personality. Anchored in the truth of fact, though seeking the truth of interpretation, biography tends to be more stable than other literary arts; and its future would appear to be a predictably steady evolution of its present trends. (P.M.K.)

Literary criticism

Construed loosely, literary criticism is the reasoned consideration of literary works and issues. It applies, as a term, to any argumentation about literature, whether or not specific works are analyzed. Plato's cautions against the risky consequences of poetic inspiration in general in his *Republic* are thus often taken as the earliest important example of literary criticism. More strictly construed, the term covers only what has been called "practical criticism," the interpretation of meaning and the judgment of quality. Criticism in this narrow sense can be distinguished not only from aesthetics (the philosophy of artistic value) but also from other matters that may concern the student of literature: biographical questions, bibliography, historical knowledge, sources and influences, and problems of method. Thus, especially in academic studies, "criticism" is often considered to be separate from "scholarship." In practice, however, this distinction often proves artificial, and even the most single-minded concentration on a text may be informed by outside knowledge, while many notable works of criticism combine discussion of texts with broad arguments about the nature of literature and the principles of assessing it. Criticism will here be taken to cover all phases of literary understanding, though the emphasis will be on the evaluation of literary works and of their authors' places in literary history. One particular aspect of literary criticism is covered in the article HISTORY, THE STUDY OF: *Textual criticism*. The following article deals with criticism largely in the context of Western literature. Its role in other literatures, along with the history of the literatures themselves, is treated in articles on the arts of various peoples, such as SOUTH ASIAN ARTS and AFRICAN ARTS.

FUNCTIONS

The functions of literary criticism vary widely, ranging from the reviewing of books as they are published to systematic theoretical discussion. Though reviews may sometimes determine whether a given book will be widely sold, many works succeed commercially despite negative reviews, and many classic works, including Herman Melville's *Moby Dick* (1851), have acquired appreciative publics long after being unfavourably reviewed and at first neglected. One of criticism's principal functions is to express the shifts in sensibility that make such revaluations possible. The minimal condition for such a new appraisal is, of course, that the original text survive. The literary critic is sometimes cast in the role of scholarly detective, unearthing, authenticating, and editing unknown manuscripts. Thus, even rarefied scholarly skills may be put to criticism's most elementary use, the bringing of literary works to a public's attention.

The variety of criticism's functions is reflected in the range of publications in which it appears. Criticism in the daily press rarely displays sustained acts of analysis and may sometimes do little more than summarize a publisher's claims for a book's interest. Weekly and biweekly magazines serve to introduce new books but are often more discriminating in their judgments, and some of these magazines, such as *The (London) Times Literary Supplement* and *The New York Review of Books*, are far from indulgent toward popular works. Sustained criticism can also be found in monthlies and quarterlies with a broad circulation, in "little magazines" for specialized audiences, and in scholarly journals and books.

Because critics often try to be lawgivers, declaring which works deserve respect and presuming to say what they are

"really" about, criticism is a perennial target of resentment. Misguided or malicious critics can discourage an author who has been feeling his way toward a new mode that offends received taste. Pedantic critics can obstruct a serious engagement with literature by deflecting attention toward inessential matters. As the French philosopher-critic Jean-Paul Sartre observed, the critic may announce that French thought is a perpetual colloquy between Pascal and Montaigne not in order to make those thinkers more alive but to make thinkers of his own time more dead. Criticism can antagonize authors even when it performs its function well. Authors who regard literature as needing no advocates or investigators are less than grateful when told that their works possess unintended meaning or are imitative or incomplete.

What such authors may tend to forget is that their works, once published, belong to them only in a legal sense. The true owner of their works is the public, which will appropriate them for its own concerns regardless of the critic. The critic's responsibility is not to the author's self-esteem but to the public and to his own standards of judgment, which are usually more exacting than the public's. Justification for his role rests on the premise that literary works are not in fact self-explanatory. A critic is socially useful to the extent that society wants, and receives, a fuller understanding of literature than it could have achieved without him. In filling this appetite, the critic whets it further, helping to create a public that cares about artistic quality. Without sensing the presence of such a public, an author may either prostitute his talent or squander it in sterile acts of defiance. In this sense, the critic is not a parasite but, potentially, someone who is responsible in part for the existence of good writing in his own time and afterward.

Although some critics believe that literature should be discussed in isolation from other matters, criticism usually seems to be openly or covertly involved with social and political debate. Since literature itself is often partisan, is always rooted to some degree in local circumstances, and has a way of calling forth affirmations of ultimate values, it is not surprising that the finest critics have never paid much attention to the alleged boundaries between criticism and other types of discourse. Especially in modern Europe, literary criticism has occupied a central place in debate about cultural and political issues. Sartre's own *What Is Literature?* (1947) is typical in its wide-ranging attempt to prescribe the literary intellectual's ideal relation to the development of his society and to literature as a manifestation of human freedom. Similarly, some prominent American critics, including Alfred Kazin, Lionel Trilling, Kenneth Burke, Philip Rahv, and Irving Howe, began as political radicals in the 1930s and sharpened their concern for literature on the dilemmas and disillusionments of that era. Trilling's influential *The Liberal Imagination* (1950) is simultaneously a collection of literary essays and an attempt to reconcile the claims of politics and art.

Such a reconciliation is bound to be tentative and problematic if the critic believes, as Trilling does, that literature possesses an independent value and a deeper faithfulness to reality than is contained in any political formula. In Marxist states, however, literature has usually been considered a means to social ends and, therefore, criticism has been cast in forthrightly partisan terms. Dialectical materialism does not necessarily turn the critic into a mere guardian of party doctrine, but it does forbid him to treat literature as a cause in itself, apart from the working class's needs as interpreted by the party. Where this utilitarian view prevails, the function of criticism is taken to be continuous with that of the state itself, namely, furtherance of the social revolution. The critic's main obligation is not to his texts but rather to the masses of people whose consciousness must be advanced in the designated direction. In periods of severe orthodoxy, the practice of literary criticism has not always been distinguishable from that of censorship.

HISTORICAL DEVELOPMENT

Antiquity. Although almost all of the criticism ever written dates from the 20th century, questions first posed

Criticism
and
scholarship

Criticism
as a target
of criticism

Criticism
in Marxist
states

by Plato and Aristotle are still of prime concern, and every critic who has attempted to justify the social value of literature has had to come to terms with the opposing argument made by Plato in *The Republic*. The poet as a man and poetry as a form of statement both seemed untrustworthy to Plato, who depicted the physical world as an imperfect copy of transcendent ideas and poetry as a mere copy of the copy. Thus, literature could only mislead the seeker of truth. Plato credited the poet with divine inspiration, but this, too, was cause for worry; a man possessed by such madness would subvert the interests of a rational polity. Poets were therefore to be banished from the hypothetical republic.

In his *Poetics*—still the most respected of all discussions of literature—Aristotle countered Plato's indictment by stressing what is normal and useful about literary art. The tragic poet is not so much divinely inspired as he is motivated by a universal human need to imitate, and what he imitates is not something like a bed (Plato's example) but a noble action. Such imitation presumably has a civilizing value for those who empathize with it. Tragedy does arouse emotions of pity and terror in its audience, but these emotions are purged in the process (*katharsis*). In this fashion Aristotle succeeded in portraying literature as satisfying and regulating human passions instead of inflaming them.

Although Plato and Aristotle are regarded as antagonists, the narrowness of their disagreement is noteworthy. Both maintain that poetry is mimetic, both treat the arousing of emotion in the perceiver, and both feel that poetry takes its justification, if any, from its service to the state. It was obvious to both men that poets wielded great power over others. Unlike many modern critics who have tried to show that poetry is more than a pastime, Aristotle had to offer reassurance that it was not socially explosive.

Aristotle's practical contribution to criticism, as opposed to his ethical defense of literature, lies in his inductive treatment of the elements and kinds of poetry. Poetic modes are identified according to their means of imitation, the actions they imitate, the manner of imitation, and its effects. These distinctions assist the critic in judging each mode according to its proper ends instead of regarding beauty as a fixed entity. The ends of tragedy, as Aristotle conceived them, are best served by the harmonious disposition of six elements: plot, character, diction, thought, spectacle, and song. Thanks to Aristotle's insight into universal aspects of audience psychology, many of his dicta have proved to be adaptable to genres developed long after his time.

Later Greek and Roman criticism offers no parallel to Aristotle's originality. Much ancient criticism, such as that of Cicero, Horace, and Quintilian in Rome, was absorbed in technical rules of exegesis and advice to aspiring rhetoricians. Horace's verse epistle *The Art of Poetry* is an urbane amplification of Aristotle's emphasis on the decorum or internal propriety of each genre, now including lyric, pastoral, satire, elegy, and epigram, as well as Aristotle's epic, tragedy, and comedy. This work was later to be prized by Neoclassicists of the 17th century not only for its rules but also for its humour, common sense, and appeal to educated taste. *On the Sublime*, by the Roman-Greek known as "Longinus," was to become influential in the 18th century but for a contrary reason: when decorum began to lose its sway encouragement could be found in Longinus for arousing elevated and ecstatic feeling in the reader. Horace and Longinus developed, respectively, the rhetorical and the affective sides of Aristotle's thought, but Longinus effectively reversed the Aristotelian concern with regulation of the passions.

Medieval period. In the Christian Middle Ages criticism suffered from the loss of nearly all the ancient critical texts and from an antipagan distrust of the literary imagination. Such Church Fathers as Tertullian, Augustine, and Jerome renewed, in churchly guise, the Platonic argument against poetry. But both the ancient gods and the surviving classics reasserted their fascination, entering medieval culture in theologically allegorized form. Encyclopaedists and textual commentators explained the supposed Christian content of pre-Christian works and the Old Testament. Although

there was no lack of rhetoricians to dictate the correct use of literary figures, no attempt was made to derive critical principles from emergent genres such as the fabliau and the chivalric romance. Criticism was in fact inhibited by the very coherence of the theologically explained universe. When nature is conceived as endlessly and purposefully symbolic of revealed truth, specifically literary problems of form and meaning are bound to be neglected. Even such an original vernacular poet of the 14th century as Dante appears to have expected his *Divine Comedy* to be interpreted according to the rules of scriptural exegesis.

The Renaissance. Renaissance criticism grew directly from the recovery of classic texts and notably from Giorgio Valla's translation of Aristotle's *Poetics* into Latin in 1498. By 1549 the *Poetics* had been rendered into Italian as well. From this period until the later part of the 18th century Aristotle was once again the most imposing presence behind literary theory. Critics looked to ancient poems and plays for insight into the permanent laws of art. The most influential of Renaissance critics was probably Lodovico Castelvetro, whose 1570 commentary on Aristotle's *Poetics* encouraged the writing of tightly structured plays by extending and codifying Aristotle's idea of the dramatic unities. It is difficult today to appreciate that this obeisance to antique models had a liberating effect; one must recall that imitation of the ancients entailed rejecting scriptural allegory and asserting the individual author's ambition to create works that would be unashamedly great and beautiful. Classicism, individualism, and national pride joined forces against literary asceticism. Thus, a group of 16th-century French writers known as the Pléiade—notably Pierre de Ronsard and Joachim du Bellay—were simultaneously classicists, poetic innovators, and advocates of a purified vernacular tongue.

The ideas of the Italian and French Renaissance were transmitted to England by Roger Ascham, George Gascoigne, Sir Philip Sidney, and others. Gascoigne's "Certain notes of Instruction" (1575), the first English manual of versification, had a considerable effect on poetic practice in the Elizabethan Age. Sidney's *Defence of Poesie* (1595) vigorously argued the poet's superiority to the philosopher and the historian on the grounds that his imagination is chained neither to lifeless abstractions nor to dull actualities. The poet "doth not only show the way, but giveth so sweet a prospect into the way, as will entice any man to enter into it." While still honouring the traditional conception of poetry's role as bestowing pleasure and instruction, Sidney's essay presages the Romantic claim that the poetic mind is a law unto itself.

Neoclassicism and its decline. The Renaissance in general could be regarded as a neoclassical period, in that ancient works were considered the surest models for modern greatness. Neoclassicism, however, usually connotes narrower attitudes that are at once literary and social: a worldly-wise tempering of enthusiasm, a fondness for proved ways, a gentlemanly sense of propriety and balance. Criticism of the 17th and 18th centuries, particularly in France, was dominated by these Horatian norms. French critics such as Pierre Corneille and Nicolas Boileau urged a strict orthodoxy regarding the dramatic unities and the requirements of each distinct genre, as if to disregard them were to lapse into barbarity. The poet was not to imagine that his genius exempted him from the established laws of craftsmanship.

Neoclassicism had a lesser impact in England, partly because English Puritanism had kept alive some of the original Christian hostility to secular art, partly because English authors were on the whole closer to plebeian taste than were the court-oriented French, and partly because of the difficult example of Shakespeare, who magnificently broke all of the rules. Not even the relatively severe classicist Ben Jonson could bring himself to deny Shakespeare's greatness, and the theme of Shakespearean genius triumphing over formal imperfections is echoed by major British critics from John Dryden and Alexander Pope through Samuel Johnson. The science of Newton and the psychology of Locke also worked subtle changes on neoclassical themes. Pope's *Essay on Criticism* (1711) is a Horatian compendium of maxims, but Pope feels obliged

Classicism, individualism, and national pride as literary forces

Horatian and Longinian tendencies

to defend the poetic rules as "Nature methodiz'd"—a portent of quite different literary inferences from Nature. Dr. Johnson, too, though he respected precedent, was above all a champion of moral sentiment and "mediocrity," the appeal to generally shared traits. His preference for forthright sincerity left him impatient with such intricate conventions as those of the pastoral elegy.

The decline of Neoclassicism is hardly surprising; literary theory had developed very little during two centuries of artistic, political, and scientific ferment. The 18th century's important new genre, the novel, drew most of its readers from a bourgeoisie that had little use for aristocratic dicta. A Longinian cult of "feeling" gradually made headway, in various European countries, against Neoclassical canons of proportion and moderation. Emphasis shifted from concern for meeting fixed criteria to the subjective state of the reader and then of the author himself. The spirit of nationalism entered criticism as a concern for the origins and growth of one's own native literature and as an esteem for such non-Aristotelian factors as "the spirit of the age." Historical consciousness produced by turns theories of literary progress and primitivistic theories affirming, as one critic put it, that "barbarous" times are the most favourable to the poetic spirit. The new recognition of strangeness and strong feeling as literary virtues yielded various fashions of taste for misty sublimity, graveyard sentiments, medievalism, Norse epics (and forgeries), Oriental tales, and the verse of plowboys. Perhaps the most eminent foes of Neoclassicism before the 19th century were Denis Diderot in France and, in Germany, Gottfried Lessing, Johann von Herder, Johann Wolfgang von Goethe, and Friedrich Schiller.

Romanticism. Romanticism, an amorphous movement that began in Germany and England at the turn of the 19th century, and somewhat later in France, Italy, and the United States, found spokesmen as diverse as Goethe and August and Friedrich von Schlegel in Germany, William Wordsworth and Samuel Taylor Coleridge in England, Madame de Staël and Victor Hugo in France, Alessandro Manzoni in Italy, and Ralph Waldo Emerson and Edgar Allan Poe in the United States. Romantics tended to regard the writing of poetry as a transcendently important activity, closely related to the creative perception of meaning in the world. The poet was credited with the godlike power that Plato had feared in him; Transcendental philosophy was, indeed, a derivative of Plato's metaphysical Idealism. In the typical view of Percy Bysshe Shelley, poetry "strips the veil of familiarity from the world, and lays bare the naked and sleeping beauty, which is the spirit of its forms."

Wordsworth's preface to *Lyrical Ballads* (1800), with its definition of poetry as the spontaneous overflow of powerful feelings and its attack on Neoclassical diction, is regarded as the opening statement of English Romanticism. In England, however, only Coleridge in his *Biographia Literaria* (1817) embraced the whole complex of Romantic doctrines emanating from Germany; the British empiricist tradition was too firmly rooted to be totally washed aside by the new metaphysics. Most of those who were later called Romantics did share an emphasis on individual passion and inspiration, a taste for symbolism and historical awareness, and a conception of art works as internally whole structures in which feelings are dialectically merged with their contraries. Romantic criticism coincided with the emergence of aesthetics as a separate branch of philosophy, and both signalled a weakening in ethical demands upon literature. The lasting achievement of Romantic theory is its recognition that artistic creations are justified, not by their promotion of virtue, but by their own coherence and intensity.

The late 19th century. The Romantic movement had been spurred not only by German philosophy but also by the universalistic and utopian hopes that accompanied the French Revolution. Some of those hopes were thwarted by political reaction, while others were blunted by industrial capitalism and the accession to power of the class that had demanded general liberty. Advocates of the literary imagination now began to think of themselves as enemies or gadflies of the newly entrenched bourgeoisie.

In some hands the idea of creative freedom dwindled to a bohemianism pitting "art for its own sake" against commerce and respectability. Aestheticism characterized both the Symbolist criticism of Charles Baudelaire in France and the self-conscious decadence of Algernon Swinburne, Walter Pater, and Oscar Wilde in England. At an opposite extreme, realistic and naturalistic views of literature as an exact record of social truth were developed by Vissarion Belinsky in Russia, Gustave Flaubert and Émile Zola in France, and William Dean Howells in the United States. Zola's program, however, was no less anti-bourgeois than that of the Symbolists; he wanted novels to document conditions so as to expose their injustice. Post-Romantic disillusion was epitomized in Britain in the criticism of Matthew Arnold, who thought of critical taste as a substitute for religion and for the unsatisfactory values embodied in every social class.

Toward the end of the 19th century, especially in Germany, England, and the United States, literary study became an academic discipline "at the doctoral level." Philology, linguistics, folklore study, and the textual principles that had been devised for biblical criticism provided curricular guidelines, while academic taste mirrored the prevailing impressionistic concern for the quality of the author's spirit. Several intellectual currents joined to make possible the writing of systematic and ambitious literary histories. Primitivism and Medievalism had awakened interest in neglected early texts; scientific Positivism encouraged a scrupulous regard for facts; and the German idea that each country's literature had sprung from a unique national consciousness provided a conceptual framework. The French critic Hippolyte Taine's *History of English Literature* (published in French, 1863-69) reflected the prevailing determinism of scientific thought; for him a work could be explained in terms of the race, milieu, and moment that produced it. For other critics of comparable stature, such as Charles Sainte-Beuve in France, Benedetto Croce in Italy, and George Saintsbury in England, historical learning only threw into relief the expressive uniqueness of each artistic temperament.

THE 20TH CENTURY

The ideal of objective research has continued to guide Anglo-American literary scholarship and criticism and has prompted work of unprecedented accuracy. Bibliographic procedures have been revolutionized; historical scholars, biographers, and historians of theory have placed criticism on a sounder basis of factuality. Important contributions to literary understanding have meanwhile been drawn from anthropology, linguistics, philosophy, and psychoanalysis. Impressionistic method has given way to systematic inquiry from which gratuitous assumptions are, if possible, excluded. Yet demands for a more ethically committed criticism have repeatedly been made, from the New Humanism of Paul Elmer More and Irving Babbitt in the United States in the 1920s, through the moralizing criticism of the Cambridge don F.R. Leavis and of the American poet Yvor Winters, to the most recent demands for "relevance."

No sharp line can be drawn between academic criticism and criticism produced by authors and men of letters. Many of the latter are now associated with universities, and the main shift of academic emphasis, from impressionism to formalism, originated outside the academy in the writings of Ezra Pound, T.S. Eliot, and T.E. Hulme, largely in London around 1910. Only subsequently did such academics as I.A. Richards and William Empson in England and John Crowe Ransom and Cleanth Brooks in the United States adapt the New Criticism to reform of the literary curriculum—in the 1940s. New Criticism has been the methodological counterpart to the strain of modernist literature characterized by allusive difficulty, paradox, and indifference or outright hostility to the democratic ethos. In certain respects the hegemony of New Criticism has been political as well as literary; and anti-Romantic insistence on irony, convention, and aesthetic distance has been accompanied by scorn for all revolutionary hopes. In Hulme conservatism and classicism were explicitly linked. Romanticism struck him as "spilt religion," a dangerous

Progress of
a cult of
"feeling"

Merging
intellectual
currents

exaggeration of human freedom. In reality, however, New Criticism owed much to Romantic theory, especially to Coleridge's idea of organic form, and some of its notable practitioners have been left of centre in their social thought.

The totality of Western criticism in the 20th century defies summary except in terms of its restless multiplicity and factionalism. Schools of literary practice, such as Imagism, Futurism, Dadaism, and Surrealism, have found no want of defenders and explicators. Ideological groupings, psychological dogmas, and philosophical trends have generated polemics and analysis, and literary materials have been taken as primary data by sociologists and historians. Literary creators themselves have continued to write illuminating commentary on their own principles and aims. In poetry, Paul Valéry, Ezra Pound, Wallace Stevens; in the theatre, George Bernard Shaw, Antonin Artaud, Bertolt Brecht; and in fiction, Marcel Proust, D.H. Lawrence, and Thomas Mann have contributed to criticism in the act of justifying their art.

Most of the issues debated in 20th-century criticism appear to be strictly empirical, even technical, in nature. By what means can the most precise and complete knowledge of a literary work be arrived at? Should its social and biographical context be studied or only the words themselves as an aesthetic structure? Should the author's avowed intention be trusted, or merely taken into account, or disregarded as irrelevant? How is conscious irony to be distinguished from mere ambivalence, or allusiveness from allegory? Which among many approaches—linguistic, generic, formal, sociological, psychoanalytic, and so forth—is best adapted to making full sense of a text? Would a synthesis of all these methods yield a total theory of literature? Such questions presuppose that literature is valuable and that objective knowledge of its workings is a desirable end. These assumptions are, indeed, so deeply buried in most critical discourse that they customarily remain hidden from critics themselves, who imagine that they are merely solving problems of intrinsic interest.

What separates modern criticism from earlier work is its catholicity of scope and method, its borrowing of procedures from the social sciences, and its unprecedented attention to detail. As literature's place in society has become more problematic and peripheral, and as humanistic education has grown into a virtual industry with a large group of professionals serving as one another's judges, criticism has evolved into a complex discipline, increasingly refined in its procedures but often lacking a sense of contact with the general social will. Major modern critics, to be sure, have not allowed their "close reading" to distract them from certain perennial questions about poetic truth, the nature of literary satisfaction, and literature's social utility, but even these matters have sometimes been cast in "value-free" empirical terms.

Recourse to scientific authority and method, then, is the outstanding trait of 20th-century criticism. The sociology of Marx, Max Weber, and Karl Mannheim, the mythological investigations of Sir James George Frazer and his followers, Edmund Husserl's phenomenology, Claude Lévi-Strauss's anthropological structuralism, and the psychological models proposed by Sigmund Freud and C.G. Jung have all found their way into criticism. The result has been not simply an abundance of technical terms and rules, but a widespread belief that literature's governing principles can be located outside literature. Jungian "archetypal" criticism, for example, regularly identifies literary power with the presence of certain themes that are alleged to inhabit the myths and beliefs of all cultures, while psychoanalytic exegetes interpret poems in exactly the manner that Freud interpreted dreams. Such procedures may encourage the critic, wisely or unwisely, to discount traditional boundaries between genres, national literatures, and levels of culture; the critical enterprise begins to seem continuous with a general study of man. The impetus toward universalism can be discerned even in those critics who are most skeptical of it, the so-called historical relativists who attempt to reconstruct each epoch's outlook and to understand works as they appeared to their first readers. Historical relativism does undermine cross-

cultural notions of beauty, but it reduces the record of any given period to data from which inferences can be systematically drawn. Here, too, in other words, uniform methodology tends to replace the intuitive connoisseurship that formerly typified the critic's sense of his role.

The debate over poetic truth may illustrate how modern discussion is beholden to extraliterary knowledge. Critics have never ceased disputing whether literature depicts the world correctly, incorrectly, or not at all, and the dispute has often had more to do with the support or condemnation of specific authors than with ascertainable facts about mimesis. Today it may be almost impossible to take a stand regarding poetic truth without also coming to terms with positivism as a total epistemology. The spectacular achievements of physical science have (with logic questioned by some) downgraded intuition and placed a premium on concrete, testable statements very different from those found in poems. Some of the most influential modern critics, notably I.A. Richards in his early works, have accepted this value order and have confined themselves to behavioristic study of how literature stimulates the reader's feelings. A work of literature, for them, is no longer something that captures an external or internal reality, but is merely a locus for psychological operations; it can only be judged as eliciting or failing to elicit a desired response.

Other critics, however, have renewed the Shelleyan and Coleridgean contention that literary experience involves a complex and profound form of knowing. In order to do so they have had to challenge Positivism in general. Such a challenge cannot be convincingly mounted within the province of criticism itself and must depend rather on the authority of antipositivist epistemologists such as Alfred North Whitehead, Ernst Cassirer, and Michael Polanyi. If it is now respectable to maintain, with Wallace Stevens and others, that the world is known through imaginative apprehensions of the sort that poetry celebrates and employs, this is attributable to developments far outside the normal competence of critics.

The pervasive influence of science is most apparent in modern criticism's passion for total explanation of the texts it brings under its microscope. Even formalist schools, which take for granted an author's freedom to shape his work according to the demands of art, treat individual lines of verse with a dogged minuteness that was previously unknown, hoping thereby to demonstrate the "organic" coherence of the poem. The spirit of explanation is also apparent in those schools that argue from the circumstances surrounding a work's origin to the work itself, leaving an implication that the former have caused the latter. The determinism is rarely as explicit or relentless as it was in Taine's scheme of race, milieu, and moment, but this may reflect the fact that causality in general is now handled with more sophistication than in Taine's day.

Whether criticism will continue to aim at empirical exactitude or will turn in some new direction cannot be readily predicted, for the empiricist ideal and its sanctuary, the university, are not themselves secure from attack. The history of criticism is one of oscillation between periods of relative advance, when the imaginative freedom of great writers prompts critics to extend their former conceptions, and periods when stringent moral and formal prescriptions are laid upon literature. In times of social upheaval criticism may more or less deliberately abandon the ideal of disinterested knowledge and be mobilized for a practical end. Revolutionary movements provide obvious instances of such redirection, whether or not they identify their pragmatic goals with the cause of science. It should be evident that the future of criticism depends on factors that lie outside criticism itself as a rationally evolving discipline. When a whole society shifts its attitudes toward pleasure, unorthodox behaviour, or the meaning of existence, criticism must follow along.

As Matthew Arnold foresaw, the waning of religious certainty has encouraged critics to invest their faith in literature, taking it as the one remaining source of value and order. This development has stimulated critical activity, yet, paradoxically, it may also be responsible in part

Multiplicity and factionalism

Literature as a form of knowing

Science as the outstanding influence

Literature as a source of order

for a growing impatience with criticism. What Arnold could not have anticipated is that the faith of some moderns would be apocalyptic and Dionysian rather than a sober and attenuated derivative of Victorian Christianity. Thought in the 20th century has yielded a strong undercurrent of anarchism which celebrates libidinous energy and self-expression at the expense of all social constraint, including that of literary form. In the critical writings of D.H. Lawrence, for example, fiction is cherished as an instrument of unconscious revelation and liberation. A widespread insistence upon prophetic and ecstatic power in literature seems at present to be undermining the complex, irony-minded formalism that has dominated modern discourse. As literary scholarship has acquired an ever-larger arsenal of weapons for attacking problems of meaning, it has met with increasing resentment from people who wish to be nourished by whatever is elemental and mysterious in literary experience.

An awareness of critical history suggests that the development is not altogether new, for criticism stands now approximately where it did in the later 18th century, when the Longinian spirit of expressiveness contested the sway of Boileau and Pope. To the extent that modern textual analysis has become what Hulme predicted, "a classical revival," it may not be welcomed by those who want a direct and intense rapport with literature. What is resisted now is not Neoclassical decorum but impersonal methodology, which is thought to deaden commitment. Such resistance may prove beneficial if it reminds critics that rationalized procedures are indeed no substitute for engagement. Excellent work continues to be written, not because a definitive method or synthesis of methods has been found, but on the contrary because the best critics still understand that criticism is an exercise of private sympathy, discrimination, and moral and cultural reflection.

(F.C.C.)

CHILDREN'S LITERATURE

Children's literature first clearly emerged as a distinct and independent form of literature in the second half of the 18th century, before which it had been at best only in an embryonic stage. During the 20th century, however, its growth has been so luxuriant as to make defensible its claim to be regarded with the respect—though perhaps not the solemnity—that is due any other recognized branch of literature.

DEFINITION OF TERMS

"Children." All potential or actual young literates, from the instant they can with joy leaf through a picture book or listen to a story read aloud, to the age of perhaps 14 or 15, may be called children. Thus "children" includes "young people." Two considerations blur the definition. Today's young teenager is an anomaly: his environment pushes him toward a precocious maturity. Thus, though he may read children's books, he also, and increasingly, reads adult books. Second, the child survives in many adults. As a result, some children's books (e.g., Lewis Carroll's *Alice in Wonderland*, A.A. Milne's *Winnie-the-Pooh*, and, at one time, Munro Leaf's *Story of Ferdinand*) are also read widely by adults.

"Literature." In the term children's literature, the more important word is literature. For the most part, the adjective imaginative is to be felt as preceding it. It comprises that vast, expanding territory recognizably staked out for a junior audience, which does not mean that it is not also intended for seniors. Adults admittedly make up part of its population: children's books are written, selected for publication, sold, bought, reviewed, and often read aloud by grown-ups. Sometimes they seem also to be written with adults in mind, as for example the popular French *Astérix* series of comics parodying history. Nevertheless, by and large there is a sovereign republic of children's literature. To it may be added five colonies or dependencies: first, "appropriated" adult books satisfying two conditions—they must generally be read by children and they must have sharply affected the course of children's literature (Daniel Defoe's *Robinson Crusoe*, Jonathan Swift's *Gulliver's Travels*, the collection of folktales by the brothers Jacob and Wilhelm Grimm, the folk-verse anthology *Des Knaben Wunderhorn* ["The Boy's Magic Horn"], edited by Achim von Arnim and Clemens Brentano, and William Blake's *Songs of Innocence*); second, books the audiences of which seem not to have been clearly conceived by their creators (or their creators may have ignored, as irrelevant, such a consideration) but that are now fixed stars in the child's literary firmament (Mark Twain's *Adventures of Huckleberry Finn*, and Charles Perrault's fairy tales); third, picture books and easy-to-read stories commonly subsumed under the label of literature but qualifying as such only by relaxed standards (though Beatrix Potter and several other writers do nonetheless qualify); fourth, first quality children's versions of adult classics (Walter de la Mare's *Stories from the Bible*, perhaps Howard Pyle's

retellings of the Robin Hood ballads and tales; finally, the domain of once oral "folk" material that children have kept alive—folktales and fairy tales; fables, sayings, riddles, charms, tongue twisters; folksongs, lullabies, hymns, carols, and other simple poetry; rhymes of the street, the playground, the nursery; and, supremely, *Mother Goose* and nonsense verse.

Five categories that are often considered children's literature are excluded from this section. The broadest of the excluded categories is that of unblushingly commercial and harmlessly transient writing, including comic books, much of which, though it may please young readers, and often for good reasons, is for the purposes of this article notable only for its sociohistorical, rather than literary, importance. Second, all books of systematic instruction are barred except those sparse examples (e.g., the work of John Amos Comenius) that illuminate the history of the subject. Third, excluded from discussion is much high literature that was not originally intended for children; from the past, Jean de La Fontaine's *Fables*, James Fenimore Cooper's *Leatherstocking* tales, Sir Walter Scott's *Ivanhoe*, Charlotte Brontë's *Jane Eyre*, Alexandre Dumas' *Three Musketeers*, Rudyard Kipling's *Kim*; from the modern period, Marjorie Kinnan Rawlings' *Yearling*, J.D. Salinger's *Catcher in the Rye*, *The Diary of Anne Frank*, Thor Heyerdahl's *Kon-Tiki*, Enid Bagnold's *National Velvet*. A fourth, rather minor, category comprises books about the young where the content but not the style or point of view is relevant (Sir James Barrie's *Sentimental Tommy*, William Golding's *Lord of the Flies*, F. Anstey's [Thomas Anstey Guthrie] *Vice Versa*). Finally, barred from central, though not all, consideration is the "nonfiction," or fact, book. Except for a handful of such books, the bright pages of which still rain influence or which possess artistic merit, this literature should be viewed from its socioeducational-commercial aspect.

THE CASE FOR A CHILDREN'S LITERATURE

Many otherwise comprehensive histories of literature slight or omit the child's reading interests. Many observers have made explicit the suspicion that children's literature, like that of detection or suspense, is "inferior." They cannot detect a sufficiently long "tradition"; distinguish an adequate number of master works; or find, to use on thoughtful critic's words, "style, sensibility, vision."

Others, holding a contrary view, assert that a tradition of two centuries is not to be ignored.

Though the case for a children's literature must primarily rest on its major writers (including a half dozen literary geniuses), it is based as well on other supports that bolster its claim to artistic stature.

Children's literature, while a tributary of the literary mainstream, offers its own identifiable, semidetached history. In part it is the issue of certain traceable social movements, of which the "discovery" of the child (see below) is the most salient. It is independent to the degree

Categories often considered children's literature

that, while it must meet many of the standards of adult literature, it has also developed aesthetic criteria of its own by which it may be judged. According to some of its finest practitioners, it is independent, too, as the only existing literary medium enabling certain things to be said that would otherwise remain unsaid or unsayable. The nature of its audience sets it apart; it is often read, especially by children younger than 12, in a manner suggesting trance, distinct from that of adult reading. Universally diffused among literate peoples, it offers a rich array of genres, types, and themes, some resembling grown-up progenitors, many peculiar to itself. Its "style, sensibility, vision" range over a spectrum wide enough to span matter-of-fact realism and tenuous mysticism.

Other measures of its maturity include an extensive body (notably in Germany, Italy, Sweden, Japan, and the United States) of commentary, scholarship, criticism, history, biography, and bibliography, along with the beginnings of an aesthetic theory or philosophy of composition. Finally, one might note its power to engender its own institutions: publishing houses, theatres, libraries, itinerant storytellers, critics, periodicals, instruction in centres of higher learning, lectureships, associations and conferences, "book weeks," collections, exhibitions, and prizes. Indeed, the current institutionalizing of children's literature on an international scale has gone so far, some feel, as to cast a shadow on the spontaneity and lack of self-consciousness that should lie at its heart.

SOME GENERAL FEATURES AND FORCES

The discovery of the child. A self-aware literature flows from a recognition of its proper subject matter. The proper subject matter of children's literature, apart from informational or didactic works, is children. More broadly, it embraces the whole content of the child's imaginative world and that of his daily environment, as well as certain ideas and sentiments characteristic of it. The population of this world is made up not only of children themselves but of animated objects, plants, even grammatical and mathematical abstractions; toys, dolls, and puppets; real, chimerical, and invented animals; miniature or magnified humans; spirits or grotesques of wood, water, air, fire, and space; supernatural and fantasy creatures; figures of fairy tale, myth, and legend; imagined familiars and doppelgänger; and grown-ups as seen through the child's eyes—whether Napoleon, Dr. Dolittle, parents, or the corner grocer. That writers did not detect this lively cosmos for two and a half millennia is one of the curiosities of literature. At any moment there has always been a numerous, physically visible, and audible company of children. Whether this sizable minority, appraised as literary raw material, could be as rewarding as the adult majority was never asked.

And so, almost to the dawn of the Industrial Revolution, children's literature remained recessive. The chief, though not the only, reason is improbably simple: the child himself, though there, was not seen—not seen, that is, as a child.

In preliterate societies he was and is viewed in the light of his social, economic, and religious relationship to the tribe or clan. Though he may be nurtured in all tenderness, he is thought of not as himself but as a pre-adult, which is but one of his many forms. Among Old Testament Jews the child's place in society replicated his father's, molded by his relation to God. So, too, in ancient Greece and Rome the child, dressed in the modified adult costume that with appropriate changes of fashion remained his fate for centuries to come, was conceived as a miniature adult. His importance lay not in himself but in what Aristotle would have called his final cause: the potential citizen-warrior. A girl child was a seedbed of future citizen-warriors. Hence classical literature either does not see the child at all or misconstrues him. Astyanax and Ascanius, as well as Medea's two children, are not persons. They are stage props. Aristophanes scorns as unworthy of dramatic treatment the children in Euripides' *Alceste*.

Throughout the Middle Ages and far into the late Renaissance the child remained, as it were, terra incognita. A sharp sense of generation gap—one of the motors of a

children's literature—scarcely existed. The family, young and old, was a kind of homogenized mix. Sometimes children were even regarded as infrahuman: for Montaigne they had "neither mental activities nor recognizable body shape." The year 1658 is a turning point. In that year a Moravian educator, Comenius, published *Orbis Sensualium Pictus* (*The Visible World in Pictures*, 1659), a teaching device that was also the first picture book for children. It embodied a novel insight: children's reading should be of a special order because children are not scaled-down adults. But the conscious, systematic, and successful exploitation of this insight was to wait for almost a century.

It is generally felt that, both as a person worthy of special regard and as an idea worthy of serious contemplation, the child began to come into his own in the second half of the 18th century. His emergence, as well as that of a literature suited to his needs, is linked to many historical forces, among them the development of Enlightenment thought (Rousseau and, before him, John Locke); the rise of the middle class; the beginnings of the emancipation of women (children's literature, unlike that for grown-ups, is in large measure a distaff product) and Romanticism, with its minor strands of the cult of the child (Wordsworth and others) and of genres making a special appeal to the young (folktales and fairy tales, myths, ballads). Yet, with all these forces working for the child, he still might not have emerged had it not been for a few unpredictable geniuses: William Blake, Edward Lear, Lewis Carroll, George MacDonald, Louisa May Alcott, Mark Twain, Collodi, Hans Christian Andersen. But, once tentatively envisaged as an independent being, a literature proper to him could also be envisaged. And so in the mid-18th century what may be defined as children's literature was at last developing.

Shifting visions of the child. Even after the child had been recognized, his literature on occasion persisted in viewing him as a diminutive adult. More characteristically, however, "realistic" (that is, nonfantastic) fiction in all countries regarded the discovered child in a mirror that provided only a partial reflection of him. There are fewer instances of attempts to present the child whole, in the round, than there are (as in Tolstoy or Joyce) attempts to represent the whole adult. Twain's Huck Finn, Erich Kästner's Emil (in *Emil and the Detectives*), Vadim Frol'ov's Sasha (in *What It's All About*), and Maria Gripe's delightful Josephine all exemplify in-the-round characterization. More frequently, however, children's literature portrays the young as types. Thus there is the brand of hell of the Puritan tradition; the moral child of Mrs. Trimmer; the well-instructed child of Madame de Genlis; the small upper class benefactor of Arnaud Berquin; the naughty child, modulated variously in Catherine Sinclair's *Holiday House* and in the books of Comtesse de Ségur, E. Nesbit, Dr. Heinrich Hoffmann (*Struwwelpeter*), and Wilhelm Busch (*Max und Moritz*); the rational child of Maria Edgeworth; the little prig of Thomas Day's *Sandford and Merton*; the little angel (Frances Hodgson Burnett's *Little Lord Fauntleroy*); the forlorn waif (Hector Malot's *Sans Famille*); the manly, outdoor child (Arthur Ransome's *Swallows and Amazons*); etc. The rationale behind these shifting visions of childhood is akin to Renaissance theories of "humours" or "the ruling passion." Progress in children's literature depended partly on abandoning this mechanical, part-for-the-whole attitude. One encouraging note in realistic children's fiction of the second half of the 20th century in all advanced countries is the appearance of a more organic view.

Slow development. A third universal feature: children's literature appears later than adult and grows more slowly. Only after the trail has been well blazed does it make use of new techniques, whether of composition or illustration. As for content, only after World War II did it exploit certain realistic themes and attitudes, turning on race, class, war, and sex, that had been part of general literature at least since the 1850s. This tardiness may be due to the child's natural conservatism.

Fourth, the tempo of development varies sharply from country to country and from region to region. It is plausible that England should create a complex children's literature, while a less-developed region (the Balkans, for

Beginnings
of
children's
literature

Realistic
themes and
attitudes

Scholarship
and
criticism

example) might not. Less clear is why the equally high cultures of France and England should be represented by unequal literatures.

The didactic versus the imaginative. The fifth, and most striking, general feature is the creative tension resulting from a constantly shifting balance between two forces: that of the pulpit-schoolroom and that of the imagination. The first force may take on many guises. It may stress received religious or moral doctrine, thus generating the Catholic children's literature of Spain or the moral tale of Georgian and early Victorian England. It may bear down less on morality than on mere good manners, propriety, or adjustment to the prevailing social code. It may emphasize nationalist or patriotic motives, as in Edmondo De Amicis' post-Risorgimento *Cuore* (*The Heart of a Child*) or much Soviet production. Or its concern may be pedagogical, the imparting of "useful" information, frequently sugarcoated in narrative or dialogue. Whatever its form, it is distinguishable from the shaping spirit of imagination, which ordinarily embodies itself in children's games and rhymes, the fairy tale, the fantasy, animal stories such as Kipling's *Jungle Books*, nonsense, nonmoral poetry, humour, or the realistic novel conceived as art rather than admonition.

Children's literature designed for entertainment rather than self-improvement, aiming at emotional expansion rather than acculturation, usually develops late. *Alice in Wonderland*, the first supreme victory of the imagination (except for *Mother Goose*), did not appear until 1865. Frequently the literature of delight has underground sources of nourishment and inspiration: oral tradition, nursery songs, and the folkish institutions of the chapbook and the penny romance.

While the didactic and the imaginative are conveniently thought of as polar, they need not always be inimical. *Little Women* and *Robinson Crusoe* are at once didactically moral and highly poetical. Nevertheless, many of the acknowledged classics in the field, from *Alice* to *The Hobbit*, incline to fantasy, which is less true of literature for grown-ups.

THE DEVELOPMENT OF CHILDREN'S LITERATURE

Criteria. Keeping these five general features of development in mind, certain criteria may now be suggested as helpful in making a gross estimate of the degree of that development within any given country. Some of these criteria are artistic. Others link with social progress, wealth, technological level, or the political structure. In what seems their order of importance, these criteria are:

1. Degree of awareness of the child's identity (see above).
2. Progress made beyond passive dependence on oral tradition, folklore, and legend.
3. Rise of a class of professional writers, as distinct from moral reformers, schoolteachers, clerics, or versatile journalists—all those who, for pedagogical, doctrinal, or pecuniary reasons turn themselves into writers for children. For example, a conscious Italian literature for young people may be said to have begun in 1776 with the Rev. Francesco Soave's moralistic "Short Stories," and largely because that literature continued to be composed largely by nonprofessionals, its record has been lacklustre. It took more than a century after the Rev. Francesco to produce a *Pinocchio*. And only in the 20th century, as typified by the outstanding work of a professional like Gianni Rodari (e.g., *Telephone Tales*), did children's literature in Italy seem to be getting into full stride.
4. Degree of independence from authoritarian controls: church, state, school system, a rigid family structure. Although this criterion might be rejected by historians of some nations, one must somehow try to explain why the Spanish, a great and imaginative people, took so long—indeed until 1952—to produce, in Sanchez-Silva, a children's writer of any notable talent.
5. Number of "classics" the influence of which transcends national boundaries.
6. Invention of new forms or genres and the exploitation of a variety of traditional ones.
7. Measure of dependence on translations.
8. Quantity of primary literature: that is, annual produc-

tion of children's books and, more to the point, of good children's books.

9. Quantity of secondary literature: richness and scope of a body of scholarship, criticism, reviewing.

10. Level of institutional development: libraries, publishing houses, associations, etc.

To these criteria some might add a vigorous tradition of illustration. But that is arguable. While Beatrix Potter's words and pictures compose an indivisible unit, it is equally true that a country may produce a magnificent school of artists (Czechoslovakia's Jirí Trnka, Ota Janeček, and others) without developing a literature of matching depth and variety.

The criteria applied: three examples. *West versus East.* The first application of such standards reveals the expected: a gap separating the achievement of the Far East from that of the West. Some Eastern literatures (New Guinea) have not advanced beyond the stage of oral tradition. Others (India, the Philippines, Ceylon, Iran) have been handicapped by language problems. Professional children's writers are rarer than in the West: according to D.R. Kalia, former director of the Delhi Public Library, "No such class exists in Hindi." In Japan, authoritarian patterns—filial piety and ancestor worship—have operated as brakes, though far less since World War II. A low economic level and inadequate technology discourage, in such countries as Burma, Sri Lanka (Ceylon), and Thailand, the origination and distribution of indigenous writing. A towering roadblock is the tendency to imitate the children's books of the West.

It is true that this vast Eastern region, considered as a whole, has produced a number of works ranking as "classics." Most advanced is Japan. Its literature for children goes back at least to the late 19th century and by 1928 was established in its own right. Japan's "discovery" of the child seems to have been made directly after World War II. In Iwaya Sazanami, Japan has its Grimm; in Ogawa Minei, perhaps its Andersen; in the contemporary Ishii Momoko, a critic and creative writer of quality; in Takeyama Michio's *Harp of Burma* (available in English), a high-quality postwar controversial novel. But, though less markedly in Japan, the basic Oriental inspiration remains fixed in folklore (also, in China and Japan, in nursery songs and rhymes), and the didactic imperative continues to act as a hobble. By most criteria the development of Eastern (as compared with Western) children's literature still appears to be sparse and tentative.

North versus south. In western Europe there is a sharp variation or unevenness, as between north and south, in the tempo of development. This basic feature was first pointed out by Paul Hazard, a French critic, in *Les Livres, les enfants et les hommes* (Eng. trans. by Marguerite Mitchell, *Books, Children and Men*, 1944; 4th ed., 1960): "In the matter of literature for children the North surpasses the South by a large margin." For Hazard, Spain had no children's literature; Italy, with its *Pinocchio* and *Cuore*, could point only to an isolated pair of works of note, and even France in order to strengthen its claims had to include northern Frenchmen: Erckmann-Chatrian, Jules Verne—and the classic Comtesse de Ségur came from Russia.

Hazard wrote in the 1920s. Since then the situation has improved, not only in his own country, but in Italy and in Portugal. Yet he is essentially correct: the south cannot match the richness of England, Scotland, Germany, and the Scandinavian countries. To reinforce his position, one might also adduce the United States, noting that the Mason-Dixon line is (though not in the field of general literature) a dividing line: the American South, even including the Uncle Remus stories, has supplied very little good children's reading. As for nursery literature, though analogous rhymes are found everywhere, especially in China, the English *Mother Goose* is unique in the claims made for it as a work of art.

Why is the north superior to the south? The first criterion of development may be illuminating. It simply restates Hazard's dictum: "For the Latins, children have never been anything but future men. The Nordics have understood better this truer truth, that men are only grown-

Japanese development

Superiority of northern literature

up children." ("Adults are obsolete children," says the American children's author "Dr. Seuss.") Hazard does not mention other factors. Historically, the south has shown greater attachment to authoritarian controls. Also, up to recent times, it has depended heavily on reworked folklore as against free invention. Besides, there is the mysterious factor of climate: it could be true that children in Latin countries mature faster and are sooner ready for adult literature. In France a special intellectual tradition, that of Cartesian logic, tends to discourage a children's literature. Clear and distinct ideas, excellent in themselves, do not seem to feed the youthful imagination.

Latin America. Again applying the chosen criteria, familiar patterns are recognizable: unevenness, as compared with the United States; belatedness—in Argentina the *cuento infantil* is hardly detectable before 1900; and especially an unbalanced polarity, with didacticism decidedly the stronger magnet. The close connection of the church with the child's family and school life has encouraged a literature stressing piety, and this at a time when the West, at least in its northern latitudes, is concerned less with the salvation than with the imagination of the child. Fantasy emerged only in the 1930s, in Brazil and in Mexico, where a Spanish exile, Antoniorrobles (pen name of Antonio Robles), continued to develop his inventive vein. And realistic writing about the actual life of the young evolved even more deliberately, being generally marked by a patriotic note. Though understandable and wholesome, this did not seem to help the cause of the imagination.

Folklore has been vigorously exploited, often by scholars of high repute. It is largely influenced by the legendry of Spain. Cuba, however, has produced interesting Afro-American tales for children; Argentina offers some indigenous folk stories and tales of gaucho life; and Central America is rich in native traditional verse enjoyed by children.

Latin American literature in general displays a special characteristic, part of its Iberian heritage: a partiality for linguistic decoration, which is unpalatable to the relatively straightforward taste of the young reader. Also the Latin-American view of the child remains tinged with a sentimentality from which many European countries and the United States had by 1914 more or less freed themselves. Thus verse for children, a medium specially cultivated in Latin America, has run to the soft, the sweet, even the lachrymose rather than to the gay, the humorous, or the sanguine—moods more congenial to the child's sensibility. This is true even of the children's verse of the Nobel Prize-winning poet Gabriela Mistral. To these two weaknesses one must add a third: the practical difficulty involved in the fact that most families cannot afford books. The absence of a powerful middle class has had a retarding effect.

Children in Latin America often complain that the authors write not for them but for their parents. They are given *lectura* ("reading matter") rather than *literatura*, which is but to say that in Latin America the admonitory note, considered so useful by church, state, and parent, continues to be sounded.

In summary, and applying the criteria: some less advanced Latin-American countries can hardly be said to have a children's literature at all. Others have produced notable writers: Brazil's José Bento Monteiro Lobato, Argentina's Ana María Berry, Colombia's Rafael Pombo, Uruguay's Horacio Quiroga. Yet the quality gap separating Latin-American children's literature from that of its northern neighbour is still wide.

HISTORICAL SKETCHES OF THE MAJOR LITERATURES

England. *Overview.* The English have often confessed a certain reluctance to say good-bye to childhood. This curious national trait, baffling to their continental neighbours, may lie at the root of their supremacy in children's literature. Yet it remains a mystery.

But, if it cannot be accounted for, it can be summed up. From the critic's vantage point, the English (as well as the Scots and the Welsh) must be credited with having originated or triumphed in more children's genres than any other country. They have excelled in the school story, two solid centuries of it, from Sarah Fielding's *The Gov-*

erness; or, The Little Female Academy (1745) to, say, C. Day Lewis' *Otterbury Incident* (1948) and including such milestones as Thomas Hughes's *Tom Brown's School Days* (1857) and Kipling's *Stalky & Co.* (1899); and the boy's adventure story, with one undebatable world masterpiece in Stevenson's *Treasure Island* (1883), plus a solid line of talented practitioners, from the Victorian Robert Ballantyne (*The Coral Island*) to the contemporary Richard Church and Leon Garfield (*Devil-in-the-Fog*); the "girls' book," often trash but possessing in Charlotte M. Yonge at least one writer of exceptional vitality; historical fiction, from Marryat's vigorous but simple *Children of the New Forest* (1847) to the even more vigorous but burnished novels of Rosemary Sutcliff; the "vacation story," in which Arthur Ransome still remains unsurpassed; the doll story, from Margaret Gatty and Richard Henry Horne to the charming fancies of Rumer Godden and the remarkable serious development of this tiny genre in Pauline Clarke's *Return of the Twelves* (1962); the realism-cum-fantasy novel, for which E. Nesbit provided a classic, and P.L. Travers a modern, formulation; high fantasy (Lewis Carroll, George MacDonald, C.S. Lewis, Alan Garner); nonsense (Carroll again, Lear, Belloc); and nursery rhymes. In Jonathan Swift's *Gulliver's Travels* and Daniel Defoe's *Robinson Crusoe*, the English furnished two archetypal narratives that have bred progeny all over the world, and in Mary Norton's Tom-Thumb-and-Gulliver-born *The Borrowers* (1952) a work of art. In Leslie Brooke (*Johnny Crow's Garden*) and Beatrix Potter (e.g., *The Tale of Peter Rabbit*) they have two geniuses of children's literature (and illustration) for very small children—probably the most difficult of all the genres. In poetry they begin at the top with William Blake and continue with Christina Rossetti, Robert Louis Stevenson, Eleanor Farjeon, Walter de la Mare, A.A. Milne, and James Reeves. In the mutation of fantasy called whimsy, Milne (*Winnie-the-Pooh*) reappears as a master. In the important field of the animal story, Kipling, with his *Jungle Books* (1894, 1895) and *Just So Stories* (1902), remains unsurpassed. Finally the English have produced a number of unclassifiable masterpieces such as Kenneth Grahame's *Wind in the Willows* (which is surely more than an animal story) and several unclassifiable writers (Mayne and Lucy Boston, for example).

The social historian, surveying the same field from a different angle, would point out that the English were the first people in history to develop not only a self-conscious, independent children's literature but also the commercial institutions capable of supporting and furthering it. He would note the striking creative swing between didacticism and delight. He would detect the sources in ballads, chap-books, nurses' rhymes, and street literature that have at critical moments prompted the imagination. What would perhaps interest him most is the way in which children's literature reflects, over more than two centuries, the child's constantly shifting position in society.

Prehistory (early Middle Ages to 1712). "Children's books did not stand out by themselves as a clear but subordinate branch of English literature until the middle of the 18th century." At least one critic has used "pre-historical" to designate all children's books published in England up to 1744, when John Newbery offered *A Little Pretty Pocket-Book*.

Before that, and as far back as the Middle Ages, children came in contact with schoolroom letters. There was the Anglo-Saxon theologian and historian the Venerable Bede, with his textbook on natural science, *De natura rerum*. There were the question-and-answer lesson books of the great English scholar Alcuin; the *Colloquy* of the English abbot Aelfric; the *Elucidarium* of the archbishop of Canterbury Anselm, often thought of as the first "encyclopaedia" for young people. Not until the mid-14th century was English (the genius of which somehow seems fitter than Latin for children's books) thought of as proper for literature. For his son "litel Lowis" Geoffrey Chaucer wrote in English the "Treatise on the Astrolabe" (1391). The English child was also afflicted, in the 15th and 16th centuries, by many "Books of Courtesy" (such as *The Babees Boke*, c. 1475), the ancestors of modern, equally ineffective manuals of conduct.

Along with these instructional works, there flourished, at least from the very early Renaissance, an unofficial or popular literature. It may not have been meant for children but—no one quite knows how—children managed to recognize it as their own. It included fables, especially those of Aesop; folk legends, such as those in the much read *Gesta Romanorum*; bestiaries, which, along with Aesop, may be ancestral to that flourishing children's genre, the animal story; romances, often clustering around King Arthur and Robin Hood; fairy tales, of which Jack the Giant Killer was the type; and nursery rhymes, probably largely orally transmitted. Perhaps the most influential underground literature consisted of the chapbooks, low-priced folded sheets containing ballads and romances (*Bevis of Southampton*, and *The Seven Champions of Christendom* [1597] were favourites), sold by wandering hawkers and peddlers. They fed the imagination of the poor, old and young, from Queen Anne's reign almost through Queen Victoria's. These native products of fancy were, in the early 18th century, reinforced by the first English translations of the classically simple French fairy tales of Charles Perrault and the more self-conscious ones of Madame D'Aulnoy.

Against this primitive literature of entertainment stands a primitive literature of didacticism stretching back to the early Middle Ages. This underwent a Puritan mutation after the Restoration. It is typified by that classic for the potentially damned child, *A Token for Children* (1671), by James Janeway. The Puritan outlook was elevated by Bunyan's *Pilgrim's Progress* (1678), which, often in simplified form, was either forced upon children or more probably actually enjoyed by them in lieu of anything better. Mrs. Overthway (in Juliana Ewing's *Mrs. Overthway's Remembrances*, 1869), recalling her childhood reading, refers to it as "that book of wondrous fascination." A softened Puritanism also reveals itself in Bunyan's *Book for Boys and Girls: or, Country Rhymes for Children* (1686), as well as the *Divine and Moral Songs for Children* by the hymn composer Isaac Watts, whose "How doth the little busy bee" still exhales a faint endearing charm.

Robinson Crusoe and Gulliver's Travels

The entire pre-1744 period is redeemed by two works of genius. Neither *Robinson Crusoe* nor *Gulliver's Travels* was meant for children. Immediately abridged and bowdlerized, they were seized upon by the prosperous young. The poorer ones, the great majority, had to wait for the beginning of the cheap reprint era. Both books fathered an immense progeny in the children's field. Defoe engendered a whole school of "Robinsonnades" in most European countries, the most famous example being Wyss's *Swiss Family Robinson* (1812-13).

On the whole, during the millennium separating Alcuin from Newbery, the child's mind was thought of, if at all, as something to be improved; his imagination as something to be shielded; his soul as something to be saved. And on the whole the child's mind, imagination, and soul resisted, persisted, and somehow, whether in a dog-eared penny history of *The Babes in the Wood* or the matchless chronicle of *Gulliver among the Lilliputians*, found its own nourishment.

From "T.W." to "Alice" (1712?-1865). Napoleon called the English a "nation of shopkeepers," and in England art may owe much to trade. Children's literature in England got its start from merchants such as Thomas Boreman, of whom little is known, and especially John Newbery, of whom a great deal more is known. Research has established that at least as early as 1730 Boreman began publishing for children (largely educational works) and that in 1742 he produced what sounds like a recreational story, *Cajanus, the Swedish Giant*. Beginnings of English children's literature might be dated from the first decade of the 18th century, when a tiny 12-page, undated book called *A Little Book for Little Children* by "T.W." appeared. It is instructional but, as the critic Percy Muir says, important as the earliest publication in English "to approach the problem from the point of view of the child rather than the adult." In sum, without detracting from the significance of Newbery, it may be said that he was merely the first great success in a field that had already undergone a certain amount of exploitation.

The elevation of the publisher-bookseller-editor Newbery (who also sold patent medicines) to the position of patron saint is an excusable piece of sentiment. Perhaps it originated with one of his back writers who doubled as a man of genius. In Chapter XVIII of *The Vicar of Wakefield* (1766), Oliver Goldsmith lauds his employer as "the philanthropic bookseller of St. Paul's Churchyard, who has written so many books for children, calling himself their friend, but who was the friend of all mankind." There is no reason to believe that Newbery was anything but an alert businessman who discovered and shrewdly exploited a new market: middle class children, or rather their parents. Nevertheless this was a creative act. In 1744 he published *A Little Pretty Pocket-Book*. Its ragbag of contents—pictures of children's games, jingles, fables, "an agreeable Letter to read from Jack the Giant Killer," plus a bonus in the form of "a Ball and a Pincushion"—are of interest only because, addressing itself single-mindedly to a child audience, it aimed primarily at diversion. Thus children's literature clearly emerged into the light of day.

The climate of Newbery's era was nevertheless more suited to a literature of didacticism than to one of diversion. John Locke's *Some Thoughts concerning Education* (1693) is often cited as an early Enlightenment emancipatory influence. But close inspection of this manual for the mental conditioning of gentlemen reveals a strong English stress on character building and practical learning. Locke thinks little of the natural youthful inclination to poetry: "It is seldom seen that anyone discovers mines of Gold or Silver in Parnassus." He does endorse, as a daring idea, the notion that a child should read for pleasure, and he recommends Aesop. But the decisive influence was not Locke's. It came from across the Channel with Rousseau's best-seller *Émile* (1762). What is positive in Rousseau—his recognition that the child should not be too soon forced into the straitjacket of adulthood—was more or less ignored. Other of his doctrines had a greater effect on children's literature. For all his talk of freedom, he provided his young *Émile* with an amiable tyrant for a teacher, severely restricting his reading to one book *Robinson Crusoe*. It was his didactic strain, exemplified in the moral French children's literature of Arnaud Berquin and Madame de Genlis, that attracted the English.

They took more easily to Rousseau's emphasis on virtuous conduct and instruction via "nature" than they did to his advocacy of the liberation of personality. Some writers, such as Thomas Day, with his long-lived *Sandford and Merton*, were avowedly Rousseauist. Others took from him what appealed to them. Sarah Kirby Trimmer, whose *Fabulous Histories* specialized in piety, opposed the presumably free-thinking Rousseau on religious grounds but was in other respects strongly influenced by him. The same is true of Anna Laetitia Barbauld, with her characteristically titled *Lessons for Children*. But Mary Martha Sherwood could hardly have sympathized with Rousseau's notion of the natural innocence of children: the author of *The History of the Fairchild Family* (1818-47) based her family chronicle on the proposition (which she later softened) that "all children are by nature evil." Of all the members of the flourishing Rousseauist or quasi-Rousseauist school of the moral tale, only one was a true writer. Maria Edgeworth may still be read.

Though the tone varies from Miss Edgeworth's often sympathetic feeling for children to Mrs. Sherwood's Savonarolan severities, one idea dominates: a special literature for the child must be manufactured in order to improve or reform him. The reigning mythology is that of reason, a mythology difficult to sell to the young.

Yet during the period from John Newbery's *Little Pretty Pocket-Book* to Lewis Carroll's *Alice in Wonderland*, children's literature also showed signs of antiselectivity. In verse there was first of all William Blake. His *Songs of Innocence* (1789) was not written for children, perhaps indeed not written for anyone. But its fresh, anti-restrictive sensibility, flowing from a deep love for the very young, decisively influenced all English verse for children. Yet the poetry the young really read or listened to at the opening of the 19th century was not Blake but *Original Poems for Infant Minds* (1804), by "Several Young Persons," in-

Newbery's publications

Influence of Rousseau on English children's literature

cluding Ann and Jane Taylor. The Taylor sisters, though adequately moral, struck a new note of sweetness, of humour, at any rate of nonpriggishness. Their "Twinkle, twinkle, little star," included in *Rhymes for the Nursery* (1806), has not only been memorized but actually liked by many generations of small children. No longer read, but in its way similarly revolutionary, was *The Butterfly's Ball and the Grasshopper's Feast* (1807), by William Roscoe, a learned member of Parliament and writer on statistics. The gay and fanciful nonsense of this rhymed satiric social skit enjoyed, despite the seeming dominance of the moral Barbaulds and Trimmers, a roaring success. Great nonsense verse, however, had to await the coming of a genius, Edward Lear, whose *Book of Nonsense* (1846) was partly the product of an emergent and not easily explainable Victorian feeling for levity and partly the issue of a fruitfully neurotic personality, finding relief for its frustrations in the noncontingent world of the absurd and the free laughter of children.

Lear's
nonsense
verse

In prose may be noted, toward the end of the period under discussion, the dawn of romantic historical fiction, with Frederick Marryat's *Children of the New Forest* (1847), a story of the English Civil War; and of the manly open-air school novel, with Thomas Hughes's *Tom Brown's School Days* (1857). A prominent milestone in the career of the "realistic" children's family novel is *Holiday House* (1839), by Catherine Sinclair, in which at last there are children who are noisy, even naughty, yet not destined for purgatory. Though Miss Sinclair's book does conclude with a standard deathbed scene, the overall atmosphere is one of gaiety. The victories in the field of children's literature may seem small, but they can be decisive. It was a small, decisive victory to have introduced in *Holiday House* an Uncle David, whose parting admonition to his nieces and nephews is: "Now children! I have only one piece of serious, important advice to give you all, so attend to me!—Never crack nuts with your teeth!"

A similar note was struck by Henry (later Sir Henry) Cole with his *Home Treasury* series, featuring traditional fairy tales, ballads, and rhymes. The fairy tale then began to come into its own, perhaps as a natural reaction to the moral tale. John Ruskin's *King of the Golden River* (1851) and William Makepeace Thackeray's "fireside pantomime" *The Rose and the Ring* (1855) were signs of a changing climate, even though the Grimm-like directness of the first is partly neutralized by Ruskin's moralistic bent and the gaiety of the second is spoiled by a laborious, parodic slyness. More important than these fairy tales, however, was the aid supplied by continental allies: the English publication in 1823–26 of the Grimms' *Fairy Tales*; in 1846 of Andersen's utterly personal fairy tales and folktales; in the '40s and '50s of other importations from the country of fancy, notably Sir George Dasent's version of the stirring *Popular Tales from the Norse* (1859), collected by Peter Christen Asbjørnsen and J.E. Moe. Though the literature of improvement continued to maintain its vigour, England was readying itself for Lewis Carroll.

Coming of age (1865–1945). In 1863 there appeared *The Water-Babies* by Charles Kingsley. In this fascinating, yet repulsive, "Fairy Tale for a Land-Baby," an unctuous cleric and a fanciful poet, uneasily inhabiting one body, collaborated. *The Water-Babies* may stand as a rough symbol of the bumpy passage from the moral tale to a lighter, airier world. Only two years later that passage was achieved in a masterpiece by an Oxford mathematical don, the Reverend Charles Lutwidge Dodgson (Lewis Carroll). *Alice's Adventures in Wonderland* improved none, delighted all. It opened what from a limited perspective seems the Golden Age of English children's literature, a literature in fair part created by Scotsmen: George MacDonald, Andrew Lang, Robert Louis Stevenson, Kenneth Grahame, James Barrie.

The age is characterized by a literary level decisively higher than that previously achieved; the creation of characters now permanent dwellers in the child's imagination (from Alice herself to Mary Poppins, and including Long John Silver, Mowgli, intelligent Mr. Toad, and—if Hugh Lofting, despite his American residence, be accepted as

English—Dr. Dolittle); the exaltation of the imagination in the work of Carroll, MacDonald, Stevenson, E. Nesbit, Grahame, Barrie, Hudson, Lofting, Travers, and the early Tolkien (*The Hobbit* [1938]); the establishment of the art fairy tale (Jean Ingelow with *Mopsa the Fairy* [1869]; Dinah Maria Mulock Craik with *The Little Lame Prince* [1875]; Mrs. Ewing with *Old Fashioned Fairy Tales* [1882]; Barrie's *Peter Pan* [1904]; and the exquisite artifices of Oscar Wilde in *The Happy Prince, and Other Tales* [1888]); the transmutation and popularization, by Andrew Lang, Joseph Jacobs, and others, of traditional fairy tales from all sources; the development of a quasi-realistic school in the fiction of Charlotte M. Yonge (*Countess Kate*); Mrs. Ewing (*Jan of the Windmill*); and Mrs. Molesworth; and, furthering this trend, a growing literary population of real, or at least more real, children (by E. Nesbit and Ransome).

It is further characterized by the rapid evolution of a dozen now-basic genres, including the school story, the historical novel, the vacation story, the "group" or "gang" novel, the boy's adventure tale, the girl's domestic novel, the animal tale, the career novel (Noel Streatfeild's *Ballet Shoes*, 1936), the work of pure whimsy (A.A. Milne's *Winnie-the-Pooh*, 1926); the solution, a brilliant one by Beatrix Potter and a charming one by L. Leslie Brooke, of the problem of creating literature for pre-readers and beginning readers; and the growth of an impressive body of children's verse: the lyric delicacy of Christina Rossetti in *Sing-Song* (1872), the accurate reflection of the child's world in Stevenson's *Child's Garden of Verses*, the satirical nonsense of Hilaire Belloc in his *The Bad Child's Book of Beasts* (1896), the incantatory, other-worldly magic of Walter de la Mare with his *Songs of Childhood* (1902) and *Peacock Pie* (1913), the fertile gay invention of Eleanor Farjeon, and the irresistible charm of Milne in *When We Were Very Young* (1924).

Finally it is characterized by the dominance in children's fiction of middle and upper middle class mores: the appearance, in the late 1930s, with Eve Garnett's *The Family from One End Street*, of stories showing a sympathetic concern with the lives of slum children; the reflection, also in the 30s, of a serious interest, influenced by modern psychology, in the structure of the child's vision of the world; the rise, efflorescence, and decline of the children's magazine: *Boy's Own Magazine* (1855–74), *Good Words for the Young* (1867–77), *Aunt Judy's Magazine* (1866–85), and—famous for its outstanding contributors—*The Boy's Own Paper* (1879–1912); the beginning, with F.J.H. Darton and other scholars, of an important critical-historical literature; institutionalization, commercialization, standardization—the popularity, for example, of the "series"; and the dominating influence of the better English work on the reading taste of American, Continental, and Oriental children.

During these 80 years a vast amount of trash and treacle was produced. What will be remembered is the work of a few dozen creative writers who applied to literature for children standards as high as those ordinarily applied to mainstream literature.

Contemporary times. If the contemporary wood cannot be seen for the trees, it is in part because the number of trees has grown so great. The profusion of English, as of children's books in general, makes judgment difficult. Livelier merchandising techniques (the spread of children's bookshops, for example), the availability of cheap paperbacks, improved library services, serious and even distinguished reviewing—these are among the post-World War II institutional trends helping to place more books in the hands of more children. Slick transformation formulas facilitate the rebirth of books in other guises: radio, television, records, films, digests, cartoon versions. Such processes may also create new child audiences, but that these readers are undergoing a literary experience is open to doubt.

Among the genres that fell in favour, the old moral tale, if not a corpse, surely became obsolescent but raised the question whether it was being replaced by a subtler form of didactic literature, preaching racial, class, and international understanding. The standard adventure story too

Post-
World War
II literature

The
Golden
Age in
English
children's
literature

seemed to be dying out, though excellent examples, such as *The Cave* (U.S. title, *Five Boys in a Cave* [1950]), by Richard Church, continued to appear. The boy's school story suffered a similar fate, despite the remarkable work of William Mayne in *A Swarm in May* (1955). Children's verse by Ian Serrailier, Ted Hughes, James Reeves, and the later Eleanor Farjeon, excellent though it was, did not speak with the master tones of a de la Mare or the precise simplicity of a Stevenson. In science fiction one would have expected more of a boom; yet nothing appeared comparable to Jules Verne.

Conversely, there was a genuine boom in fact books: biographical series, manuals of all sorts, popularized history, junior encyclopaedias. Preschool and easy-to-read beginners' books, often magnificently produced, multiplied. So did specially prepared decoys for the reluctant reader. After the discovery of the child came that of the postchild: conscientiously composed teen-age and "young adult" novels were issued in quantity, though the quality still left something to be desired. A 19th-century phenomenon—experimentation in the juvenile field by those who normally write for grown-ups—took on a second life after World War II. Naomi Mitchison, Richard Church, P.H. Newby, Richard Graves, Eric Linklater, Norman Collins, Roy Fuller, C. Day Lewis, and Ian Fleming, with his headlong pop extravaganza *Chitty Chitty Bang Bang* (1964), come to mind.

A post-World War II stress on building bridges of understanding was reflected both in an increase in translations and in the publication of books, whether fiction or non-fiction, dealing responsibly and un sentimentally with the sufferings of a war-wounded world. One example among many was Serrailier's *Silver Sword* (1958), recounting the trans-European adventures that befell four Polish children after the German occupation. *The Silver Sword* was a specialized instance of a general trend toward the interpretation for children of a postwar world of social incoherence, race and class conflict, urban poverty, and even mental pathology. Such novels as John Rowe Townsend's *Gumble's Yard* (1961); *Widdershins Crescent* (1965); *Pirate's Island* (1968); Eve Garnett's *Further Adventures of the Family from One End Street* (1956); and Leila Berg's *Box for Benny* (1958) represented a new realistic school, restrained in England, less so in the United States, but manifest in the children's literature of much of the world. It failed to produce a masterpiece, perhaps because the form of the realistic novel must be moderately distorted to make it suitable for children.

In two fields, however, English postwar children's literature set new records. These were the historical novel and that cloudy area comprising fantasy, freshly wrought myth, and indeed any fiction not rooted in the here and now.

There was fair reason to consider Rosemary Sutcliff not only the finest writer of historical fiction for children but quite unconditionally among the best historical novelists using English. A sound scholar and beautiful stylist, she made few concessions to the presumably simple child's mind and enlarged junior historical fiction with a long series of powerful novels about England's remote past, especially that dim period stretching from pre-Roman times to the coming of Christianity. Among her best works are *The Eagle of the Ninth* (1954), *The Shield Ring* (1956), *The Silver Branch* (1957), *The Lantern Bearers* (1959), and especially *Warrior Scarlet* (1958).

Not as finished in style, but bolder in the interpretation of history in terms "reflecting the changed values of the age," was the pioneering Geoffrey Trease. He also produced excellent work in other juvenile fields. Typical of his highest energies is the exciting *Hills of Varna* (1948), a story of the Italian Renaissance in which Erasmus and the great printer Aldus Manutius figure prominently. Henry Treece, whose gifts were directed to depicting violent action and vigorous, barbaric characters, produced a memorable series of Viking novels of which *Swords from the North* (1967) is typical.

This new English school, stressing conscientious scholarship, realism, honesty, social awareness, and general disdain for mere swash and buckle, produced work that completely eclipsed the rusty tradition of Marryat and

George Alfred Henty. Some of its foremost representatives were Cynthia Harnett, Serrailier, Barbara Leonie Picard, Ronald Welch (pseudonym of Ronald O. Felton), C. Walter Hodges, Hester Burton, Mary Ray, Naomi Mitchison, and K.M. Peyton, whose "Flambards" series is a kind of Edwardian historical family chronicle. Leon Garfield, though not working with historical characters, created strange picaresque tales that gave children a thrilling, often chilling insight into the 18th-century England of Smollett and Fielding.

In the realm of imagination England not only retained but enhanced its supremacy with such classics as *Tom's Midnight Garden* (1958), by Ann Philippa Pearce, a haunting, perfectly constructed story in which the present and Victoria's age blend into one. There is the equally haunting Green Knowe series, by Lucy M. Boston, the first of which, *The Children of Greene Knowe*, appeared when the author was 62. The impingement of a world of legend and ancient, unsleeping magic upon the real world is the basic theme of the remarkable novels of Alan Garner. Complex, melodramatic, stronger in action than in characterization, they appeal to imaginative, "literary" children. Garner's rather nightmarish narrative *The Owl Service* (1967) is perhaps the most subtle.

Finally there is a trio of masters, each the architect of a complete secondary world. The vast Middle Earth trilogy *The Lord of the Rings* (1954–55), by the Anglo-Saxon and Middle English language scholar J.R.R. Tolkien, was not written with children in mind. But they have made it their own. It reworks many of the motives of traditional romance and fantasy, including the Quest, but is essentially a structure, conceivably but not inevitably allegorical, of sheer invention on a staggering scale. It is also a sociocultural phenomenon, selling 3,000,000 copies in nine languages and functioning, for a certain class of American teenagers, as a semisacred cult object.

Tolkien's fellow scholar, C.S. Lewis, created his own otherworld of Narnia. It is more derivative than Tolkien's (he owes something, for example, to Nesbit), more clearly Christian-allegorical, more carefully adapted to the tastes of children. Though uneven, the seven volumes of the cycle, published through the years 1950 to 1956, are exciting, often humorous, inventive, and, in the final scenes of *The Last Battle*, deeply moving.

The third of these classic secondary worlds is in a sense not a creation of fantasy. The four volumes (1952–61) about the Borrowers, with their brief pendant, *Poor Stainless* (1971), ask the reader to accept only a single impossibility, that in a quiet country house, under the grandfather clock, live the tiny Clock family: Pod, Homily, and their daughter Arrietty. All that follows from this premise is logical, precisely pictured, and carries absolute conviction. Many critics believe that this miniature world so lovingly, so patiently fashioned by Mary Norton will last as long as those located at the bottom of the rabbit hole and through the looking glass.

United States. *Overview.* Compared with England, the United States has fewer peaks. In *Huckleberry Finn*, of course, it possesses a world masterpiece matched in the children's literature of no other country. *Little Women*, revolutionary in its day, radiates a century later a special warmth and may still be the most beloved "family story" ever written. Though *The Wonderful Wizard of Oz* has been recklessly compared with *Alice*, it lacks Carroll's brilliance, subtlety, and humour. Nonetheless, its story and characters apparently carry, like *Pinocchio*, an enduring, near-universal appeal for children. To these older titles might be added *Stuart Little* (1945) and *Charlotte's Web* (1952), by E.B. White, two completely original works that appear to have become classics. To this brief list of high points few can be added, though, on the level just below the top, the United States bears comparison with England and therefore any other country.

The "law" of belated development applies in a special way. From Jamestown to the end of the Civil War, American children's literature virtually depended on currents in England. In the adult field Cooper and Washington Irving may stand for a true declaration of independence. But it was not until the 1860s and '70s, with Mary

Teenage
novels

Novels of
Tolkien
and C.S.
Lewis

Historical
fiction for
children

American
classics

Mapes Dodge's *Hans Brinker*, Louisa May Alcott's *Little Women*, Lucretia Hale's *Peterkin Papers*, Mark Twain's *Tom Sawyer*, and *St. Nicholas* magazine, that children's literature finally severed its attachment to the mother country. In the marketplace, however, a uniquely American note was sounded much earlier, the first of the Peter Parley series of Samuel Goodrich having appeared in 1827.

In certain important fields, the United States pioneered. These include everyday-life books for younger readers; the non-class-based small-town story such as *The Moffats* by Eleanor Estes; the Americanized fairy tale and folktale such as *Uncle Remus* (1880), not originally meant for children, and Carl Sandburg's *Rootabaga Stories* (1922); beginners' books such as Dr. Seuss's *The Cat in the Hat* (1957); and the "new realism." One might maintain that American children's literature, particularly that since World War II, is bolder, more experimental, more willing to try and fail, than England's. Moreover, it set new standards of institutionalization, "packaging," merchandising, and publicity, as well as mere production, especially of fact books and "subject series."

Prehistory (1646?-1865). The prehistoric annals are short and simple. Dominated by England, native creativity—to refer only to books with even the thinnest claim to literary quality—amounted to little. The Puritan view of the unredeemable child obtained almost into the era of Andrew Jackson. Jonathan Edwards put it neatly: unrepentant children were "young vipers and infinitely more hateful than vipers." More moderate notions also existed. Imported English ballads and tales, even a few "shockers," were enjoyed by the young vipers. But in general, from John Cotton's *Spiritual Milk for Boston Babes* (1646) through the Civil War, the admonitory and exemplary tract and the schoolmaster's pointer prevailed. Occasionally there is the cheerful note of non-improvement, as in Clement Moore's "Visit from St. Nicholas" (1823), sounding against the successful lesson-cum-moral tales of Peter Parley (Goodrich) and the didactic "Rollo" series of Jacob Abbott. The latter's *Franconia Stories* (1850-53), however, showing traces of Rousseau and Johann Pestalozzi, is the remote ancestor of those wholesome, humorous pictures of small-town child life in which American writers excelled after World War I. Affectionately based on the author's own memories, they occasionally reveal children rather than improvable miniatures of men.

The children's magazines of the early 19th century did their best to amuse as well as instruct the young. Sara Josepha Hale's "Mary Had a Little Lamb" appeared in *The Juvenile Miscellany* (1826-34). The atmosphere was further lightened by *Grandfather's Chair* (1841) and its sequels, retellings of stories from New England history by Nathaniel Hawthorne. These were followed in 1852-53 by his redactions, rather unacceptable today, of Greek legends in *The Wonder Book for Girls and Boys* and *Tanglewood Tales for Girls and Boys*. Hawthorne's death date (1864) coincided roughly with a qualified subsidence of the literature of the didactic.

Peaks and plateaus (1865-1940). During the period from the close of the Civil War to the turn of the century an Americanized white, Anglo-Saxon, Protestant, Victorian gentility dominated as the official, though not necessarily real, culture. At first glance such a climate hardly seems to favour the growth of a children's literature. But counterforces were at work: a vigorous upsurge of interest, influenced by European thinkers, in the education and nurture of children; the dying-out of the old Puritanism; and the accumulation of enough national history to stimulate the imagination. To these forces must be added the appearance in Louisa May Alcott of a minor genius and in Samuel Clemens (Mark Twain) of a major one.

American materialism (and also its optimism) expressed itself in the success myth of Horatio Alger, while a softened didacticism, further modified by a mild talent for lively narrative, was reflected in the 116 novels of Oliver Optic (William Taylor Adams). But a quartet of books appearing from 1865 to 1880—heralded a happier day. These were Mary Mapes Dodge's *Hans Brinker, or the Silver Skates* (1865), which for all its Sunday-school tone, revealed to American children an interesting foreign culture and told

a story that still has charm; Louisa May Alcott's *Little Women* (1868; vol. ii, 1869; and its March family sequels), which lives by virtue of the imaginative power that comes from childhood truly and vividly recalled; Lucretia Hale's *Peterkin Papers* (1880), just as funny today as a century ago, perfect nonsense produced in a non-nonsensical era; and Thomas Bailey Aldrich's *Story of a Bad Boy* (1870). This, it is often forgotten, preceded *Tom Sawyer* by seven years, offered a model for many later stories of small-town bad boys, and is a fair example of the second-class classic. But it took *Tom Sawyer* and *Huckleberry Finn* to change the course of American writing and give the first deeply felt vision of boyhood in juvenile literature.

To these names should be added Frank Stockton (whose *Ting-a-Ling Tales* [1870] showed the possibilities inherent in the invented fairy tale) and especially the writer-illustrator Howard Pyle. His reworkings of legend (*The Merry Adventures of Robin Hood*, 1883; the King Arthur stories, 1903-1910, and his novels of the Middle Ages [*Otto of the Silver Hand*, 1888; and *Men of Iron*, 1892]) exemplify perfectly the romantic feeling of his time, as does the picture of Shakespeare's England drawn by John Bennett in *Master Skylark* (1897).

The sentimentality that is sometimes an unconscious compensatory gesture in a time of ruthless materialism expressed itself in the idyllic *Poems of Childhood* (1896), by Eugene Field, and the rural dialect *Rhymes of Childhood* (1891), by James Whitcomb Riley. These poems can hardly speak to the children of the second half of the 20th century. But it is not clear that the same is true of the equally sentimental novels of Frances Hodgson Burnett. It is easy to smile over *Little Lord Fauntleroy* (1886) or her later and superior novels, *A Little Princess* (1905) and *The Secret Garden* (1911). Back of the absurd sentimentality, however, lies an extraordinary narrative skill, as well as an ability to satisfy the perennial desire felt by children at a certain age for life to arrange itself as a fairy tale.

The development of a junior literature from 1865 to about 1920 is ascribable less to published books than to two remarkable children's magazines: *The Youth's Companion* (1827-1929, when it merged with *The American Boy*) and the relatively nondidactic *St. Nicholas* magazine (1873-1939), which exerted a powerful influence on its exclusively respectable child readers. (It is surely needless to point out that up to the 1960s children's literature has been by and for the middle class.) These magazines published the best material they could get, from England as well as the United States. For all their gentility, standards, including that of illustration, were high. The contributors' names in many cases became part of the canon of world literature. To the children of the last quarter of the 19th and first quarter of the 20th century, the periodical delivery of these magazines presumably meant something that film and television cannot mean to today's children. The magazines were not "media." They were friends.

Appropriately the new century opened with a novelty: a successful American fairy tale. *The Wonderful Wizard of Oz* (1900) is vulnerable to attacks on its prose style, incarnating mediocrity. But there is something in it, for all its doctrinaire moralism, that lends it permanent appeal: a prairie freshness, a joy in sheer invention, the simple, satisfying characterization of Dorothy and her three old, lovable companions. Several of the sequels—but only those bearing L. Frank Baum's name—are not greatly inferior.

The century underwent for the next two decades a rather baffling decline. Some institutional progress was made in library development, professional education, and the reviewing of children's books. Much useful work was also accomplished in the field of fairy-tale and folktale collections. But original literature did not flourish. There were Pyle and Mrs. Burnett and the topflight nonsense verses of Laura E. Richards, whose collected rhymes in *Tirra Litra* (1932) will almost bear comparison with those of Edward Lear. Less memorable are the works of Lucy Fitch Perkins, Joseph Altsheler, Ralph Henry Barbour, Kate Douglas Wiggin, Eliza Orne White, and the two Burgesses—Thornton and Gclett. During these decades, de la Mare, Miss Potter, Kipling, Barrie, Grahame, and E. Nesbit were at work in England.

Contributions of Howard Pyle

Development of a junior literature

Dominance of Victorian gentility

During the period between world wars new trails were blazed in nonfiction with van Loon's *Story of Mankind* and V.M. Hillyer's *Child's History of the World* (1922). The *Here and Now Story Book*, by Lucy Sprague Mitchell, published in the 1920s, was the first real example of the "direct experience" school of writing, but it is more properly part of the chronicle of pedagogy than of literature. The small child was far better served by a dozen talented writer-illustrators, such as Wanda Gág, with her classic *Millions of Cats* (1928) and other delightful books; and Ludwig Bemelmans, with *Madeline* (1939) and its sequels. Other distinguished names in the important and growing picture-book field were Marjorie Flack, Hardie Gramatky, James Daugherty, the d'Aulaires, and Virginia Lee Burton.

In the field of comic verse and pictures for children of almost all ages, Dr. Seuss (Theodore Geisel), starting with *And to Think That I Saw It on Mulberry Street* (1937), continued to lead, turning out so many books that one tended to take him for granted. His talent is of a very high order.

Literature of the 1920s and '30s

The 1920s and '30s produced many well-written historical novels, striking a new note of authority and realism, such as *Drums* (1925, transformed in 1928 into a boy's book with N.C. Wyeth's illustrations), by James Boyd, and *The Trumpeter of Krakow* (1928), by Eric Kelly. The "junior novel" came to the fore in the following decade, together with an increase in books about foreign lands, minority groups, and a boom in elaborate picture books. Children's verse was well served by such able practitioners as Dorothy Aldis and Rosemary and Stephen Vincent Benét, with their stirring, hearty ballad-like poems collected in *A Book of Americans* (1933). But the only verse comparable to that of Stevenson or de la Mare was the exquisite *Under the Tree* (1922), by the novelist Elizabeth Madox Roberts, a treasure that should never be forgotten.

At least three other writers produced work of high and entirely original quality. Two of them—Florence and Richard Atwater—worked as a pair. Their isolated effort, *Mr. Popper's Penguins* (1938), will last as a masterpiece of deadpan humour that few children or adults can resist. The third writer is Laura Ingalls Wilder. Her *Little House* books, nine in all, started in 1932 with *The Little House in the Big Woods*. The entire series, painting an unforgettable picture of pioneer life, is a masterpiece of sensitive recollection and clean, effortless prose.

Work of quality was contributed during these two lively decades by authors too numerous to list. Among the best of them are Will James, with his horse story *Smoky* (1926); Rachel Field, whose *Hitty* (1929) is one of the best doll stories in the language; Elizabeth Coatsworth, with her fine New England tale *Away Goes Sally* (1934); and the well-loved story of a New York tomboy in the 1890s, *Roller Skates* (1936), by the famous oral storyteller Ruth Sawyer.

Contemporary times. Since the 1930s the quality and weight of American children's literature were sharply affected by the business of publishing, as well as by the social pressures to which children, like adults, were subjected. Intensified commercialization and broad-front expansion had some good effects and some bad ones as well.

For any book of interest to adults, publishers constructed a corresponding one scaled to child size. The practice of automatic miniaturization stimulated a pullulation of fact books—termed by an unsympathetic observer "the information trap"—marked by a flood of subject series and simplified technology. Paperbacks and cheap reprints of juvenile favourites enlarged the youthful reading public, just as the multiplication of translations widened its horizon. More science fiction was published, a field in which the stories of Robert Heinlein and *A Wrinkle in Time* (1962), by Madeleine L'Engle, stood out. An increase was also noticeable in books for the disadvantaged child and in work of increasingly high quality by and for blacks. In the early 1950s, children's book clubs flourished, though they appeared to be on the wane little more than a decade later. Simple narration using "scientifically determined vocabulary" also seemed to decrease in popularity. There was a marked tendency to orient titles, fiction and nonfiction,

to the requirements of the school curriculum. Another trend was toward collaborative "international" publishing. This had the double effect of cutting colour-plate costs and promoting blandness, since it was important that no country's readers be offended or surprised by anything in text or illustration. Still another alteration took place in the conventional notion of age and grade levels. Teenagers reached out for adult books; younger children read junior novels.

The most striking development was the growth of the "realists," most of them as earnest as Maria Edgeworth, a few of them lighter fingered, with a fringe of far-outers. The latter were fairly represented by Nat Hentoff in *Jazz Country* (1965), for example, and Maria Wojciechowska in *The Rotten Years* (1971). Teenage fiction as well as nonfiction dealt mercilessly with ethnic exploitation, poverty, broken homes, desertion, unemployment, adult hypocrisy, drug addiction, sex (including homosexuality), and death. A whole new "problem" literature became available, with no sure proof that it was warmly welcomed. The aesthetic dilemmas posed by this literature are still to be faced and resolved. The new social realist story often had the look of an updated moral tale: the dire consequences of nondiligence were replaced by those of pot smoking.

Growth of "realist" fiction

Nevertheless such original works as *Harriet the Spy* (1964) and *The Long Secret* (1965), by Louise Fitzhugh, showed how a writer adequately equipped with humour and understanding could incorporate into books for 11-year-olds subjects—even menstruation—ordinarily reserved for adult fiction. Similarly trailblazing were the semidocumentary novels of Joseph Krumgold: . . . *And Now Miguel* (1953), *Onion John* (1958), and *Henry 3* (1967), the last about a boy with an I.Q. of 154 trying to get along in a society antagonistic to brains. The candid suburban studies of E.L. Konigsburg introduced a new sophistication. Her 1968 Newbery Medal winner, *From the Mixed-Up Files of Mrs. Basil E. Frankweiler*, was original in its tone and humour.

As for the more traditional genres, a cheering number of high-quality titles rose above the plain of mediocrity. The nonfantastic animal story *Lassie Come Home* (1940), by Eric Knight, survived adaptation to film and television. In the convention of the talking animal, authentic work was produced by Ben Lucien Burman, with his wonderful "Catfish Bend" tales (1952–67). The American-style, wholesome, humorous family story was more than competently developed by Eleanor Estes, with her "Moffat" series (1941–43) and *Ginger Pye* (1951); Elizabeth Enright, with her Melendy family (1941–44); and Robert McCloskey, with *Homer Price* (1943)—to name only three unflinching popular writers. Text-and-picture books for the very young posed an obdurate challenge: to create literature out of absolutely simple materials. That challenge, first successfully met by Beatrix Potter, attracted Americans. The modern period produced many enchanting examples of this tricky genre: *The Happy Lion* (1954) and its sequels, the joint work of the writer Louise Fatio and her artist husband, Roger Duvoisin; the "Little Bear" books, words by Else Holmelund Minarik, pictures by Maurice Sendak; and several zany tours de force by Dr. Seuss, including his one-syllable revolution *The Cat in the Hat* (1957). The picture books of Sendak, perhaps one of the few original geniuses in his restricted field, were assailed by many adults as frightening or abnormal. The children did not seem to mind.

Fiction about foreign lands boasted at least one modern American master in Meindert De Jong, whose most sensitive work was drawn from recollections of his Dutch early childhood. A Hans Christian Andersen and Newbery winner, he is best savoured in *The Wheel on the School* (1954), and especially in the intuitive *Journey from Peppermint Street* (1968). The historical novel fared less well in America than in England. *Johnny Tremain* (1943), by Esther Forbes, a beautifully written, richly detailed story of the Revolution, stood out as one of the few high points, as did *The Innocent Wayfaring* (1943), a tale of Chaucer's England by the equally scholarly Marchette Chute. Poetry for children had at least two talented representatives. One was the eminent poet-critic John Ciardi, the other

Fiction about foreign lands

David McCord, a veteran maker of nonsense and acrobat of language.

In fantasy, the farcical note was struck with agreeable preposterousness by Oliver Butterworth in *The Enormous Egg* (1956) and *The Trouble with Jenny's Ear* (1960). The prolific writer-illustrator William Pène Du Bois has given children nothing more uproariously delightful than *The Twenty-one Balloons* (1947), merging some of the appeals of Jules Verne with those of Samuel Butler's *Erewhon* and adding a sly humour all his own. Two renowned *New Yorker* writers, James Thurber and E.B. White, developed into successful fantasists, Thurber with an elaborate series of ambiguous literary fairy tales such as *The Thirteen Clocks*, White with his pair of animal stories *Stuart Little* and *Charlotte's Web* that for their humanity and uninsistent humour stand alone. The vein of "high fantasy" of the more traditional variety, involving magic and the construction of a legendary secondary world, was represented by the five highly praised volumes of the Prydain cycle (1964-68) by Newbery Medal winner Lloyd Alexander.

Two other works of pure imagination gave the 1960s some claim to special notice. The first was *The Phantom Tollbooth* (1961) by Norton Juster, a fantasy about a boy "who didn't know what to do with himself." Not entirely unjustly, it has been compared to *Alice*. The second received less attention but is more remarkable: *The Mouse and His Child* (1969), by Russell Hoban, who had been a successful writer of gentle tales for small children. But here was a different affair altogether: a flawlessly written, densely plotted story with quiet philosophical overtones. It involved a clockwork mouse, his attached son, and an unforgettable assortment of terribly real, humanized animals. Like *Alice* and *The Borrowers*—indeed like all major children's literature—it offered as much to the grown-up as to the young reader. With this moving, intellectually demanding fantasy the decade ended on a satisfactory note.

Germany and Austria. A. Merget's *Geschichte der deutschen Jugendliteratur* ("History of German Children's Literature") appeared in 1867, some years before the Germans had much children's literature to consider, a demonstration of Teutonic thoroughness. By two criteria—degree of awareness of the child's identity and level of institutional development—Germany leads the world. It has built a vast structure of history, criticism, analysis, and controversy devoted to a subject the chief property of which would appear to be its charm rather than its obscurity. One estimate has it that in West Germany alone there are over 300 associations dedicated to the study and promotion of children's literature. Such conscientiousness, nowhere else matched, such a serious desire to relate the child's reading to his nurture, education, and *Weltanschauung*, has an admirable aspect. But by attaching juvenile books too closely to the theory and demands of pedagogy, it may have constricted a marked native genius.

The dominant historical influences roughly coincide with those that have affected German mainstream literature, though, as expected, they were exerted more slowly. The Reformation, stressing the Bible, the catechism, and the hymnbook, bent the literature of childhood toward the didactic, the monitory, and the pious. The Enlightenment, however, did something to help toward the identification of the child as an independent being. With this insight are associated the educational theories of J.B. Basedow, J.F. Herbart, and Friedrich Froebel. One fruit of the movement was *Robinson der Jüngere* (1779; "The Young Robinson"), by Joachim Heinrich Campe, who adapted Defoe along Rousseauist lines, his eye sharply fixed on what he considered to be the natural interests of the child. Interchapters of useful moral conversations between the author and his pupils were a feature of the book. Campe's widespread activities on behalf of children, though less commercially motivated, recall Newbery's.

Rationalism, piety, and the German partiality for disciplined conduct were modified by the influence of two crucial works, not intended for children but soon taken over by them. Both are part of the Romantic movement that swept Germany and much of the Continent during the early 19th century. *Des Knaben Wunderhorn* (1805-08; "The Youth's Magic Horn"), a collection of old Ger-

man songs and folk verse, included many children's songs, or songs that were so denominated by the editors, Achim von Arnim and Clemens Brentano. The effect of the book was to retrieve for Germany much of its rich folk heritage, to promote a new emotional sensibility, and to draw attention to the link, as the Romantics thought, binding folk feeling to the child's vision of the world. *Des Knaben Wunderhorn* became a part of German childhood, as La Fontaine's *Fables* in France and *Mother Goose* in England had become a part of growing up in those countries. It helped inspire several excellent writers of verse for children: A.H. Hoffmann von Fallersleben; August Kopisch; the writer-illustrator Count Franz Pocci, the first German to write nonsense verse for the young; F.W. Güll; and later poets such as Paula and Richard Dehmel.

Just as *Des Knaben Wunderhorn* became a source of poetry, so the epochal folktale collection of the brothers Grimm helped to develop a school of prose fairy-tale writers. Not all of these Romantics wrote with children in mind. But some of the simplest of their tales have become part of the German child's inheritance. In today's presumably practical era, they are once more in favour. Among these masters of the "art" *Märchen* are E.T.A. Hoffmann; C.M. Brentano; Ludwig Tieck; de la Motte Fouqué, author of *Undine*; and Wilhelm Hauff, whose talents are most nearly adapted to the tastes of children.

Two curious half-geniuses of comic verse and illustration wrote and drew for the hitherto neglected small child. *Struwwelpeter* ("Shock-headed Peter"), by the premature surrealist Heinrich Hoffmann, aroused cries of glee in children across the continent. Wilhelm Busch created the slapstick buffoonery of Max and Moritz, the ancestors of the Katzenjammer Kids and indeed of many aspects of the comic strip.

The second half of the 19th century saw an increase in commercialized sentimentality and sensation and a corresponding decline in quality. The bogus Indian and Wild West tales of Karl May stand out luridly in the history of German children's literature. Up to about 1940, 7,500,000 of his books had been sold to German readers alone. (Emilio Salgari in Italy, G.A. Henty in England, and "Ned Buntline" in the United States, who were contemporaneously satisfying the same hunger for the suspenseful, did not approach May's talent for fabrication without the slightest root in reality.)

It may have been May and others like him who roused an educator, Heinrich Wolgast, to publish in 1896 his explosive *Das Elend unserer Jugendliteratur* ("The Sad State of Our Children's Literature"). The event was an important one. It advanced for the first time the express thesis that "Creative children's literature must be a work of art"; Wolgast resolutely decried nationalistic and didactic deformations. He precipitated a controversy the echoes of which are still audible. On the whole his somewhat excessive zeal had a wholesome effect.

Two post-Wolgast poets of childhood worthy of mention are Christian Morgenstern, whose macabre, pre-Dada poetry for adults later came into vogue, and the lesser-gifted Joachim Ringelnatz. The nondidactic note they sounded in modern times was strengthened by a whole school of children's poets. No other country produced work in this difficult field superior to the finest verse of the multi-talented James Krüss, and especially Josef Guggenmos, whose lyric simplicity at times recalls Blake. Guggenmos also has to his credit a translation of *A Child's Garden of Verses*, in itself an original work of art.

Between the world wars, prose showed few high points and, after the advent of Hitler, many low ones. *Der Kampf der Tertia* (1927; "The Third-Form Struggle"), by Wilhelm Speyer, was Germany's excellent contribution to the genre of the school story. Erich Kästner's *Emil and the Detectives* (1929) ranked not only as a work of art, presenting city boys with humour and sympathy, but as an immediate classic in an entirely new field, the juvenile detective story (Mark Twain's awkward *Tom Sawyer, Detective* [1896] may be ignored). Kästner, the dean of German writers for children, won an international audience with a long series of stories of which the thesis-fable *Die Konferenz der Tiere* (1949; Eng. trans. *The Animals'*

German
fairy tales

Historical
influences

Conference, 1949) is perhaps the funniest as well as the most serious.

Post-World War II literature

Post-World War II literature, recovering from the Nazi blight, was strong in several fields. In realistic fantasy there is *Vevi* (1955), by the Austrian Erica Lillegg, an extraordinary tale of split personality, odd, exciting, even profound. Michael Ende's *Jim Knopf und Lucas der Lokomotivführer* (1961; Eng. trans., *Jim Button and Luke the Engine Driver*, 1963) has more than a touch of *Oz*; and both Kästner and Krüss have made agreeable additions to the realm of fantasy.

In the domain of the historical novel, Hans Baumann is a distinguished name. Lacking the narrative craft of Miss Sutcliff, whose story lines are always clean and clear, he matched her as a scholar and mounted scenes of great intensity in such novels as *Die Barke der Brüder* (1956; Eng. trans., *The Barque of the Brothers*, 1958) and especially *Steppensöhne* (1954; Eng. trans., *Sons of the Steppe*, 1958), a tale about two grandsons of Genghis Khan. His narrative history of some exciting archaeological discoveries, *Die Höhlen der grossen Jäger* (1953; Eng. trans., *The Caves of the Great Hunters*, 1954; rev. ed., 1962), is a minor classic. Mention should be made of Fritz Mühlenweg, a veteran of the Sven Hedin expedition of 1928–32 to Inner Mongolia and the author of *Grosser-Tiger und Kompass-Berg* (1950; Eng. trans., *Big Tiger and Christian*, 1952). A long, richly coloured narrative of a journey made by two boys, Chinese and European, through the Gobi Desert, it should stand as one of the finest adventure stories of the postwar years.

One general conclusion regarding West German children's literature after 1945 was that the native genius, which had been impeded by pedagogical theory and nationalist dogma, again appeared to be in free flow.

In East Germany, production was conditioned by the association with the Soviet Union, and it appeared to be recapitulating the developments in children's literature that had occurred in the Soviet Union after 1917. Socialist Realism was the basic food offered to the literary appetites of young East Germans.

Scandinavia. Sweden. Scandinavia, but especially Sweden, inevitably suggests a question as to why a group of small, sparsely populated countries ranks directly after England and the United States for the variety, vigour, and even genius of its children's literature. Hazard's north-south theory describes; it does not explain. A few possible factors may be listed: the inspiration of the master Andersen—yet he does not seem greatly to have inspired his homeland; the appearance in 1900 of the Swedish Ellen Key's two-volume *Barnets århundrede* (Eng. trans., *The Century of the Child*, 1909), pivotal in the history of the discovery that children really exist; a general modern atmosphere of social enlightenment; welfare statism tempered by regard for the individual; a school and library system, notably in Sweden, of extraordinary humanity and efficiency; perhaps even the long, lively career of the Stockholm Children's Theatre, a centre of creative activity. Yet the mystery persists. Since the first half of the 19th century, Scandinavia produced Andersen, Zacharias Topelius, Jorgen Moe, Henrik Wergeland, Helena Nyblom, Selma Lagerlöf, Elsa Beskow, Astrid Lindgren, Tove Jansson, Maria Gripe, Anna Lisa Warnlöf, Lennart Hellsing, Karin Anckarsvärd, Inger Sandberg, plus a school of critics and historians second only to that of Germany, plus many talented illustrators.

Children's literature in Sweden for centuries reflected that of Germany, of which Sweden was a cultural province during the Reformation and even through the Enlightenment period. The historian Göte Klingberg traced some kind of religious-instructive reading for children back to 1600. There is a record, though the manuscripts have vanished, of children's plays produced at the country manors during the 1700s and into the following century. The tradition of children's theatre has always been stronger in Sweden than elsewhere in Europe.

True Swedish national literature

A true native literature is usually dated from 1751–53, when the tutor Count Carl Tessin wrote his "Old Man's Letters to a Young Prince" (Gustav III), in which instruction was tempered by the first fairy tales written for

Swedish children. The German influence, however, persisted until about the middle of the 19th century, when Fredrika Bremer, traveller and feminist, tried to stimulate the work of indigenous children's writers. The dominant influence of the Finnish-born but basically Swedish Topelius, of Hans Christian Andersen, and of the romantic spirit in general was felt at this time. Later in the century two followers of Andersen—Helena Nyblom and Anna Wahlenberg—enriched the tradition of the fairy tale. The former's *Sagokrans* (1903; Eng. trans., *The Witch of the Woods*, 1968), preserves a rare charm.

The great landmark, however, is Miss Lagerlöf's world classic *Nils Holgerssons underbara resa genom Sverige*, 2 vol. (1906–07; Eng. trans., *The Wonderful Adventures of Nils*, 1907; *Further Adventures of Nils*, 1911). Written (at the request of the state ministry of education) as a school geography, it is the rare example of an officially commissioned book that turned out to be a work of art. *Nils*, for all its burden of instruction, is a fantasy. At the same time, a realistic breakthrough was achieved by Laura Fitinghoff, whose historical novel about the famine of the 1860s, *Barnen från Frostmojället* (1907; Eng. trans., *Children of the Moor*, 1927), ranks as a classic.

According to the historian Eva von Zweigbergk, didacticism ("diligence, obedience, and moderation") obtained up to the 1920s, though she also views the period 1890–1915 as Sweden's Golden Age. It included not only *Nils* but the emergence of a school of creators of picture books for small children headed by Elsa Beskow, whose work in pictures and text, extending over the years from 1897 to 1952, was decisive in its influence. This pre-modern period also saw many good writers for grown-ups devoting their talents to juvenile fiction. The sailing story *Målarpirater* (1911; "The Pirates of Lake Mälaren"), by the novelist Sigfrid Siwertz, is a still-remembered example.

The period from 1940 on has called forth a bewildering array of talented writers and artist-writers. In the field of humour and nonsense there are Åke Holmberg, with his parodic Ture Sventon detective series; the outstanding poet Lennart Hellsing, with *Daniel Doppsko* (1959); Astrid Lindgren, successful in a half dozen genres but perhaps best known as the creator of the supergirl Pippi Longstocking; Gösta Knutsson, with her well-liked *Pelle Svanslös* (1939; Eng. trans., *The Adventures of the Cat Who Had No Tail*). The psychological realistic novel, delving deeply into the inner lives of children, has been developed by Maria Gripe, whose *Hugo and Josephine* trilogy may become classic; Gunnel Linde's *Tacka vet jag Skorstensgränd* (1959; Eng. trans., *Chimney-Top Lane*, 1965); and Anna Lisa Warnlöf, writing under the pseudonym of "Claque," whose two series about Pella and Fredrika show an intuitive understanding of lonely and misunderstood children.

Harry Kullman and Martha Sandwall-Bergström are among the few Swedish writers who have used working class industrial backgrounds successfully. Kullman is also a historical novelist. The prolific Edith Unnerstad has written charming family stories, with a touch of fantasy, as has Karin Anckarsvärd, whose *Doktors pojke* (1963; Eng. trans., *Doctor's Boy*, 1965) is a quietly moving tale of small-town life in the horse-and-buggy days. The Sandbergs, Inger and Lasse, have advanced the Beskow tradition in a series of lovely picture books. Fantasy has been well served by Lindgren, Edith Unnerstad, Holmberg, Hellsing, and others. Children's poetry is a lively contemporary art, one distinguished poet being Britt G. Hallqvist.

By most criteria of development the Swedes rank high among those creating a children's literature that is both broad and deep.

Norway. Norway cannot boast a genius of worldwide fame. But, beginning with the 1830s when a new literary language, based on spoken Norwegian, was forged, Norway has possessed an identifiable children's literature. From 1837 to 1844 Asbjørnsen and Moe, the Grimms of Norway, published their remarkable collection of folk stories, and thus created not only a literary base on which the future could build but a needed sense of national identity. Moe also wrote specifically for children. His poems are part of Norwegian childhood, and his nature fantasy *I bronden*

og i tjernet ("In the Well and the Lake," 1851) made Viggo and his little sister Beate familiar for more than a century. Equally enduring are the fairy tales and children's verse of Norway's greatest poet Henrik Wergeland.

Golden
Age of
Norwegian
literature

The Norwegian critic Jo Tenfjord believes that the 30 years from 1890 to 1920 represented a golden age. With this period are associated Dikken Zwillgmeyer, author of the "Inger Johanne" series about a small-town little girl; Barbra Ring, creator of the popular "Peik" stories and of a play *The Princess and the Fiddler*, which was produced yearly at the National Theatre in Oslo; Gabriel Scott; and the fairy-tale writer Johan Falkberget.

Among the more prominent and well-loved moderns are Halvor Floden, whose most famous work, centred on a gypsy waif, is *Gjenta fra lands vejen* ("The Girl from the Road"); the nonsense versifier Zinken Hopp; the poet Jan-Magnus Bruheim, three of whose collections have won state prizes; Finn Havrevold, whose toughminded boys' teenage novel *Han Var Min Ven* became available in English translation as *Undertow* in 1968, and who also wrote successfully for girls; Leif Hamre, specializing in air force adventures; the prolific, widely translated Aimée Sommerfelt, whose works range from "puberty novels" to faraway stories set in Mexico City and northern India; Thorbjørn Egner, who is the author of, among other books, a tiny droll fantasy, *Karius and Baktus* (1958; Eng. trans. 1962), which will actually persuade small children to brush their teeth; and Alf Prøysen, creator of Mrs. Pepperpot, a delightful little old lady who never knows when she is going to shrink to peppercorn size. Fantasy of this kind seems less characteristic of contemporary Norway than does the realistic novel, especially that designed for older children.

Denmark. Without Hans Christian Andersen, Danish children's literature might have fared better. It is not that his countrymen deify him, as much as it is that the outside world does. Indeed, because modernized versions of his tales do not exist, his now rather antiquated Danish tends to outmode him. Yet his gigantic shadow must have intimidated his literary descendants, just as Dante and Cervantes intimidated theirs. Doubtless other forces also account for the sparseness and relative conventionality of Danish children's literature.

The earliest books were written for the children of the nobility. Not till the passage of the Education Act of 1814 did the poorer ones have access to any suitable reading matter, and this, obedient to the prevailing European fashion, was dour in tone. The climate, of course, relaxed when Andersen appeared with his phenomenal series, still the finest of their kind, of invented or reworked fantastic tales. In 1884 H.V. Kaalund published a picture book of "Fables for Children" based on the popular verse narratives (1833) of a Thuringian pastor, Wilhelm Hey. Three years later an unidentified Danish humorist added three cautionary tales to a translation of six *Struwwelpeter* stories. Though it does not seem to have appeared as a picture book until 1900, Christian Winther in 1830 wrote a pleasing trifle, with an unusual fantastic touch, called "Flugten til Amerika" ("Flight to America"). It is still ranked as a classic. Such are some of the 19th-century oases.

Denmark's general tendency has been to over-rely on translations or adaptations, drawn especially from its neighbour Germany. As against this, it can point to an excellent original tradition of nursery and nonsense rhymes. The first such collection, made as early as 1843, stimulated not only Andersen but such other 19th-century figures as Johan Krohn, whose "Peter's Christmas" remains a standard seasonal delight. The tradition is relayed to the 20th century by Halfdan Rasmussen, whose collected *Bjørnerim* ("Verse for Children") won the 1964 Danish Children's Book Prize, and Ib Spang Olsen, with his nonsense picture book *The Boy in the Moon* (1962). As for the complementary prose tradition of fireside tales, Denmark had to wait (Andersen was artist, not scholar) for its Grimm until 1884, when a collection made by Svend Grundtvig, the son of N.F.S. Grundtvig, a great bishop-educator, was posthumously published.

As compared with other Scandinavian countries, post-World War II developments lagged. Picture books exhibited much more originality than did teenage literature.

Jytte Lyngbirk's girls' novels, notably the love story "Two Days in November," however, are well reputed, as are the realistic fictions, laid against an industrial background, of Tove Ditlevsen. Perhaps Denmark's boldest original talent is Anne Holm, who aroused healthy controversy with her (to some) shocking narrative of a displaced boy's journey to Denmark, the novel *David* (1963; Eng. trans., *North to Freedom*, 1965).

Some informed observers ascribe Denmark's only moderate performance to domination by the teaching profession, to the lingering influence of conventional didacticism, and to the lack of the economic-social forces that stimulate professional writers. As late as 1966 the Minister of Culture commented on the scarcity of Danish juvenile authors, and this at a time when the rest of Scandinavia was, as it remained, in the full flood of the modern movement.

Finland. Although its language and people are not of European origin, Finland is loosely conceived as part of the Scandinavian bloc. Only since December 6, 1917, has it been formally independent. During much of its history Swedish was the language of the educated class. Thus its two outstanding premodern children's writers, the father figure Zacharias Topelius and Anni Swan, wrote their fairy tales and folktales primarily for a Swedish-reading audience. Their works however were promptly translated into Finnish and became part of the native heritage. The same is true of the contemporary Tove Jansson, 1966 Andersen Medal winner, whose series of novels about the fantastic self-contained world of Moomintrolls, though less successful with English-reading children, enchants young readers throughout northern and central Europe.

The labours of Topelius in the children's field and of Elias Lönnrot (compiler of the great Finnish epic-miscellany the *Kalevala*, 1835) in the field of national folklore constituted the soil from which Finnish children's literature was eventually to derive nutriment. But that literature emerged as an identifiable whole only after World War I. It is largely folktales rooted. Indeed this small country became an international focus of folklore research. One student has said that it probably possesses the largest number of folktales in existence, some 30,000 of them. In the early 1960s a fairy tale competition yielded 795 manuscripts, a phenomenal statistic in view of Finland's sparse population.

Finland, despite the fact that its language tends to limit its audience, is part of the main current of children's literature, even though only Jansson has won anything like an international reputation. Two children's poets, Aila Meriluoto and Kirsi Kunnas, have achieved renown.

France. *Overview.* The French themselves are not happy with their record. Writing in the late 1940s, critic Jean de Trigon, in *Histoire de la littérature enfantine, de ma Mère l'Oye au Roi Babar* (Paris, Librairie Hachette, 1950) said: "The French have created little children's literature. They have received more than they have given, but they have assimilated, adapted, transformed. The two are not the same thing, for one must love childhood in general if one is to please children other than one's own." In 1923 Marie-Thérèse Latzarus tolled the passing bell in *La littérature enfantine en France dans la seconde moitié du XIX^e siècle* (Paris: Les Presses Universitaires de France): "Children's literature, more's the pity, is dying." And in 1937, in their introduction to *Beaux livres, belles histoires*, the compilers Marguerite Gruny and Mathilde Leriche wrote: "Children's literature in France is still poor, despite the earnest efforts of the last decade."

Surely Trigon was too severe. Even more surely Mille Latzarus has proved a false Cassandra. As for the compilers, the very decade they scorned saw at least three magnificent achievements. The first was Jean de Brunhoff's. Equally talented as author and artist, in 1931 he gave the world that enlightened monarch Babar the Elephant, one of the dozen or so immortal characters in children's literature. The next year saw the start of Paul Faucher's admirable Père Castor series, imaginatively conceived, beautifully designed educational picture books for the very young—not literature, perhaps, but historically comparable to Comenius. Finally, in 1934 appeared the first of Marcel Aymé's miraculous stories about two little girls and the talking animals whose adventures they shared. These

Criticism
of French
children's
literature

Early
nursery
and
nonsense
rhymes

grave-comic fantasies were later collected as *Les Contes du chat perché* (1939; Eng. trans., *The Wonderful Farm*, 1951; *Return to the Wonderful Farm*, 1954), and, along with de Brunhoff and Faucher, were enough to make the decade great.

But there are no other decades to match it. There does exist a disproportion between French literary genius as a whole and the children's literature it has been able to produce. The explanation is uncertain. Mme Le Prince de Beaumont, an adventurous 18th-century lady who wrote over 70 volumes for the young, thought that children's stories should be pervaded by "the spirit of geometry." It is possible that the blame for France's showing might in part be laid on a persistent Cartesian spirit, reinforced by rationalist and positivist philosophies. The Cartesian does not readily surrender to fancy, especially of the more wayward variety. And so, even counting Charles Perrault, the later Charles Nodier, and the contemporary Simone Ratel and Maurice Vauthier, a dearth of first-rate fairy tales may be noted. Cartesians would tend to be weak also in children's verse, in nonsense of any sort, in humour (despite Babar), even in the more imaginative kind of historical novel exemplified by Hans Baumann in Germany and Rosemary Sutcliff in England. Perhaps French children's literature has been restrained by a Catholicism or by a Protestantism that continued to insist on the edifying when mainstream literature had already freed itself from explicit moralism. It may not even be true, as Trigon thinks, that the French have fruitfully assimilated the children's literature of foreign countries. *Alice* has more or less bewildered them; *Huckleberry Finn* has never been digested. The child's cause was not much aided by the triumph of a post-Napoleonic bourgeois cast of thought—or by the wave of post-1871 nationalism.

It is a complicated problem. But perhaps the heart of it lies in the value the French set on maturity. For them childhood at times has seemed less a normal human condition than a handicap. The children themselves have often seemed to feel the pressure, which may account for the fact that they absorb French adult books precociously. The French came much later than did many other countries to the discovery of the child as a figure worthy of the most sensitive understanding; that is what makes Père Castor so important. One is not surprised to note the comparatively recent date (1931) of a study by Aimé Dupuy, translatable as *The Child: A New Character in the French Novel*.

History. If one skips Jean de La Fontaine, whose *Fables* (1668; 1678–79; and 1693), though read by the young, were not meant for them, French children's literature from one point of view begins with the classic fairy tales of Charles Perrault. These were probably intended for the salon rather than the nursery, but their narrative speed and lucidity commended them at once to children. The fairy tales of his contemporary Mme d'Aulnoy, like many others produced in the late 17th and early 18th centuries, are hardly the real thing. With a Watteau-like charm, they taste of the court, as does the *Télémaque* of François Fénelon, a fictionalized lecture on education.

Rousseau, as has been noted, did make a difference. *Émile* at least drew attention to what education might be. But the effect on children's literature was not truly liberating. His disciple, Mme de Genlis, set a stern face against make-believe of any sort; all marvels must be explained rationally. Her stories taught children more than they wanted to know, a circumstance that endeared her to a certain type of parent. Sainte-Beuve, to be fair, called her "the most gracious and gallant of pedagogues." One of her qualities, priggishness, was energetically developed by Arnaud Berquin in his *Ami des enfants*. Berquin created the French equivalent of the concurrent English bourgeois morality. In effect, he unconsciously manufactured an adult literature for the young, loading the dice in favour of the values held by parents to be proper for children. Yet one must beware of judging Berquin or his equally moralistic successor Jean-Nicolas Bouilly by today's standards. Children accepted them because they were the best reading available; and Anatole France's tribute in *Le Petit Pierre* (1918) shows that they must have exerted some charm.

The didactic strain, if less marked than in England or

Germany, persisted throughout most of the 19th century. To it, Mme de Ségur, in her enormously popular novels, added sentimentality, class snobbery, but also some liveliness and occasional fidelity to child nature. Her "Sophie" series (1850s and 60s), frowned on by modern critics, is still loved by obstinate little French girls. *Sans Famille* (1878), by Hector Malot, a minor classic of the "unhappy child" school, also continues to be read and is indeed a well-told story. But the century's real writer of genius is of course Jules Verne, whose first book, *Un Voyage en ballon*, was originally published in 1851 in a children's magazine, *Le Musée des Familles*.

The period was lively enough. Production was vast. Children's magazines flourished, particularly the remarkable *Magasin d'Éducation et de Récréation*, brilliantly edited by Jules Hetzel. Writers of the stature of George Sand, Alphonse Daudet, and Alexandre Dumas père were not too proud to write for children. Much worthy, though transient, work was produced along with a mass of mediocrity, as was the case also in England and the United States. But on the whole, as the century drew to a close, French children might have been better served, even though one critic sees the apogee as occurring between 1860 and 1900.

From the turn of the century to the close of World War II, a number of superior works were produced. The books of de Brunhoff and Faucher have already been cited. A remarkable picture of prehistoric life by J.-H. Rosny (pseudonym of J.-H.-H. Boex) appeared in 1911 and has proved so durable that in 1967 an English translation, *The Quest for Fire*, appeared. *Patapoufs et filififers*, by André Maurois, a gentle satire on war, has lasted (Eng. trans. *Pattypuffs and Thinifers*, 1948; reissued 1968). His fantastic *Le Pays des 36,000 volontés* is almost as popular. The famous dramatist Charles Vildrac has done much to advance the cause of French children's literature. Two pleasant stories of his, remotely descended from *Robinson Crusoe*. *L'Isle rose* and its sequel *La Colonie*, appeared in the 1920s and 1930s. In 1951 his now-classic comic animal tale *Les Lunettes du lion* won immediate success (Eng. trans., *The Lion's Eyeglasses*, 1969). On a high literary level, not accessible to all children, was *Le Petit Prince* (1943, both French and English, *The Little Prince*) by the famous aviator-author Antoine de Saint-Exupéry. The very vagueness of this mystical parable has lent it a certain magnetism. Finally, it is necessary to mention a field in which the French proved incomparable: the comic strip combining action and satire, conceived on a plane of considerable sophistication. Hergé's *Tintin* started in the 1930s and sold over 25,000,000 copies. Also successful was the later and even more unconventional *Astérix* series.

Production after 1945 so multiplied that to single out names is bound to involve some injustice. A few, however, by reason either of the originality of their talent or the scope of their achievement, stand out. One is Maurice Druon, whose *Tistou of the Green Fingers* (1957; Eng. trans. 1958), a kind of children's *Candide*, demonstrated how the moral tale, given sufficient sensitivity and humour, can be transmuted into art. Perhaps the most original temperament was that of Henri Bosco, author of four eerie, haunting Provencal novels about the boy Pascalet and his strange involvements with a gypsy companion, a fox, and a dog in a shifting, legend-shrouded natural world. It may be that time will rate these books, like those of the English writer Walter de la Mare, among the finest of their kind. Bosco's *L'Enfant et la rivière* (1955; Eng. trans., *The Boy and the River*, 1956), *Le Renard dans l'île* (1956; Eng. trans., *The Fox in the Island*, 1958), and *Barboche* (1957; Eng. trans. 1959) are notable.

Sound, realistic novels, almost free of excess moralism, were written by at least a dozen reputable authors. Among them Colette Vivier (*The House of the Four Winds*), Paul-Jacques Bonzon (*The Orphans of Simitira*), and Étienne Cattin (*Night Express*) were distinguished. The domain of the imaginative tale was well represented by Maurice Vauthier, especially by his *Écoute, petit loup*. Among those noted for their prolific output as well as the high level of their art two names emerged. One is Paul Berna, who has worked in half a dozen genres, including detective stories and science fiction. His *Cheval sans tête* (1955) was pub-

Early 20th-century French literature

Perrault's fairy tales

lished in England as *A Hundred Million Francs* and in the United States as *The Horse Without a Head* and was made into a successful Disney film. A "gang" story, using a hard, unemotional tone that recalls Simenon, it may be the best of its kind since *Emil and the Detectives*.

The death of René Guillot removed a deeply conscientious and responsible artist. Guillot, though probably not of the first rank, was not far below it. He left more than 50 widely translated novels for the young and about 10 nonfiction works. For his entire body of work, he received in 1964 the Andersen Prize. His finest achievements in the adventure novel, based on his experiences in Africa, include *The White Shadow* (1948) and *Riders of the Wind* (1953).

Children's verse has at least one delightful practitioner in Pierre Gamarra. His *Mandarine et le Mandarin* contains Fontainesque fables of notable drollery and high technical skill. The Belgian author Maurice Carême also has some repute as a children's poet. In summary, contemporary French activity seems a bit lacking in colour and versatility. But one solid achievement must be registered: the 19th century's legacy was decisively rejected, and at last a natural child prevailed in the imaginative work of the best French contemporaries.

Russia/Soviet Union. Here history breaks cleanly into two periods: pre-1917 and post-1917. In pre-Revolutionary Russia may be observed a most dramatic illustration of the disproportion that may exist between a children's and a mainstream literature. Beyond question the latter is one of the greatest of the modern world. But Russia's pre-1917 children's literature is anemic. It does include the fables of Ivan Krylov; a great treasury of Russian folktales (*skazki*) assembled by A.N. Afanasyev; the epic tales (*byliny*) sung or told to children; the classic by Pyotr Yrshov, *Konyok gorbunok* (1834; English adaption by Ireen Wicker, *The Little Hunchback Horse*, 1942); and other stories and poems enjoyed by young Russians but not originally designed for them. To this folk material should be added the McGuffeyish moral tales that Tolstoy wrote for a series of graded readers. There is also the poet-translator Vasily Zhukovsky, praised by the respected critic Vissarion Belinsky as one of the few poets of the century, part of whose work was dedicated to children.

On the whole, however, pre-Revolutionary Russia could make only a few feeble gestures toward the creation of an independent children's literature. The submerged peasantry relied on the fireside tale teller. The middle class, while far stronger than is generally recognized, was in no position to stimulate or support a literature for its children. The privileged class looked to the West: the children read Mme de Genlis. Thus it came about that the child was recognized later in Russia than in other parts of western Europe. The critic and children's writer Korney Chukovsky speaks of the "indifference" with which "early childhood was regarded in the past." He then points out that attitudes have changed, so that now the child is "an adored hero."

The Revolution was the watershed. After 1917 Soviet children's literature developed more or less in accord with the necessities of the state. This is not to say that it became identical with Soviet propaganda. Indeed one of the finest teenage novels, Vadim Frolov's *Chto k chemu* (Eng. trans., *What It's All About*, 1965), is quite untouched by dogma of any kind. Soviet children's literature, and especially its vast body of popularized science and technology for the young, however, was in general governed by the ideals of socialist realism, the idolization of the "new Soviet man" (as in the widely read works of Boris Zhitkov and Arkady Gaydar), the exaltation of the machine over the irresponsible furniture of fairyland, and especially a revised version of the pre-18th-century miniature adult view of the child: he now had become a potential Soviet citizen and architect of the Communist future.

Juvenile fiction and biography naturally tended to cue themselves into the crucial episodes of Soviet history. But the theory underlying this basically nationalist literature (suggesting similar developments in Italy and England in the latter half of the 19th century) is by no means clear-cut. The most influential thinker was Maksim Gorky,

who during the 1920s called for "creative fantasy," for children's stories "which make out of the human being, instead of a will-less creature or an indifferent workman, a free and active artist, creator of a new culture." He asked for books that would encourage the child to become "a knight of the spirit." Gorky's essays are a curious, endearing mixture of Marxist doctrine (with a utopian slant) and quite standard Western humanistic ideas. It is in Korney Chukovsky's remarkable book *Malenkiye deti* (1925) or *Ot dvukh do pyati* (Eng. trans., *From Two to Five*, 1963), however, that the opposition of two familiar forces, entertainment and instruction, can be sensed most clearly. The tension is typically expressed in Chukovsky's account of the Soviet war over the fairy tale, the opposition to which reached its high point in the 1920s and '30s. "We propose," wrote one journalist in a Moscow magazine in 1924, "to replace the unrealistic folktales and fantasies with simple realistic stories taken from the world of reality and from nature." Chukovsky, himself a writer full of humour and invention, opposed this view, as had Gorky before him.

Though rich in folklore drawn from its many peoples and languages, Soviet culture remained weak in the realm of fantasy. A fairy play such as Marshak's *Krugly god* (Eng. trans., *The Month Brothers*, 1967) seems (at least in English) fatally heavy-handed. That Soviet children's literature was vigorous, varied, and motivated by a genuine concern for the child is undoubted. However, there certainly existed no Soviet "Narnia" series, a Soviet *Borrowers*.

It is not difficult to see that contemporary children's literature in Russia is lively, copious, and probably enjoyed. It is much more difficult for those who have no Russian to judge its value. Occasionally in translation one will come across something as superb as the beautiful nature and animal tales in *Arcturus the Hunting Hound and Other Stories* (1968) by Yury Kazakov. But one can only record, without judging, the vast production of such popular children's writers as Samuil Marshak, Sergey Mikhalkov, Lev Kassil, and N. Nosov. Especially notable is the popularity of poetry, whether it be the work of such past generation writers as Vladimir Mayakovsky or that of the contemporary Agniya Barto. Apparently Russian children read poetry with more passion and understanding than do English-speaking children. The mind of the Russian child is carefully looked after. He is provided with books, often beautifully illustrated, which many Western countries may find hard to match. "Demand from them as much as possible, respect them as much as possible," says Anton Makarenko, the theorist of children's literature. (C.Fa.)

BIBLIOGRAPHY

The scope of literature. *General works:* KENNETH BURKE, *The Philosophy of Literary Form*, 2nd ed. (1967); I.A. RICHARDS, *Practical Criticism: A Study of Literary Judgment* (1929, reprinted 1968) and *Principles of Literary Criticism* (1924, reprinted 1961); GEORGE SAINTSBURY, *A History of Criticism and Literary Taste in Europe from the Earliest Texts to the Present Day*, 3 vol. (1900-04, reprinted 1961); NOWELL C. SMITH (ed.), *Literary Criticism* (1905); KONSTANTIN KOLENDA, *Philosophy in Literature* (1982).

Ancient to modern: *A Translation of the Latin Works of Dame Alighieri* (1904), see Letter X to Can Grande; CHARLES SEARS BALDWIN, *Ancient Rhetoric and Poetic, Interpreted from Representative Works* (1924), *Medieval Rhetoric and Poetic to 1400, Interpreted from Representative Works* (1928, reprinted 1971), and *Renaissance Literary Theory and Practice* (1939); EDWARD H. BLAKENEY (ed.), *Horace on the Art of Poetry* (1928); CECIL MAURICE BOWRA, *Primitive Song* (1962); S.I.L. BUTCHER, *Aristotle's Theory of Poetry and Fine Art*, with a critical text and translation of the *Poetics*, 4th ed. (1907; reprinted with corrections, 1932); INGRAM BYWATER, *Aristotle on the Art of Poetry* (1909), reprinted in *Aristotle's Poetics and Longinus on the Sublime*, ed. by CHARLES SEARS BALDWIN (1930); PIERRE CORNEILLE, *Oeuvres*, 3 vol. (1862), containing the "Discourse on Dramatic Poetry" and "Discourse on the Three Unities"; ALBERT S. COOK (ed.), *The Art of Poetry* (1892, reprinted 1926), containing a translation of Horace's *Art of Poetry*; J.D. DENNISTON, *Greek Literary Criticism* (1924, reprinted 1971), translations, beginning with Aristophanes; *Dryden's Essays on the Drama*, ed. by WILLIAM STRUNK (1898); ALLAN H. GILBERT (ed.), *Literary Criticism: Plao to Dryden* (1940); EDMUND D.

JONES (ed.), *English Critical Essays (Sixteenth, Seventeenth, and Eighteenth Centuries)* (1922); Ben Jonson: *Timber, Discoveries, and Conversations with Drummond of Hawthornden* in his *Works*, ed. by C.H. HEREORD and PERCY SIMPSON, 11 vol. (1925-52); LONGINUS, *On the Sublime*, Greek text with an English translation by W. HAMILTON EYEE ("Loeb Classical Library," 1927); Plato, trans. by LANE COOPER (1938), contains the *Phaedrus*, the *Symposium*, the *Ion*, the *Gorgias*, and parts of the *Republic* and the *Laws*; GEORGE PUTTENHAM, *The Art of English Poesie*, ed. by GLADYS D. WILLCOCK and ALICE WALKER (1936); PAUL RADIN, *Primitive Man as Philosopher* (1927); *Sidney's Apologie for Poetrie*, ed. by J. CHURTON COLLINS (1907); G. GREGORY SMITH (ed.), *Elizabethan Critical Essays*, 2 vol. (1904); GAY WILSON ALLEN and H.H. CLARK (eds.), *Literary Criticism: Pope to Croce* (1941 and 1962); MATTHEW ARNOLD, *Essays in Criticism*, 2 vol., First and Second Series complete (1902), and *Essays in Criticism*, with an introduction by E.J. O'BRIEN, Third Series (1910); EDWIN BERRY BURGUM (ed.), *The New Criticism: An Anthology of Modern Aesthetics and Literary Criticism* (1930); SAMUEL TAYLOR COLFRIDGE, *Biographia Literaria, or Biographical Sketches of My Literary Life and Opinions*, 2 vol., reprinted from the original plates (1907); BENEDETTO CROCE, *The Defence of Poetry, Variations on the Theme of Shelley*, trans. by E.E. CARRITT (1933); T.S. FLIOT, *Selected Essays, 1917-1932* (1932); *Hazlitt on English Literature: An Introduction to the Appreciation of Literature*, ed. by JACOB ZEITLIN (1913, reprinted 1970); E.R. HUGHES (trans.), *The Art of Letters: Lu Chi's "Wen Fu," A.D. 302* "Bollingen Series XXIX" (1951); THOMAS ERNEST HULME, *Speculations: Essays on Humanism and the Philosophy of Art*, ed. by HERBERT READ (1924); JAMES GIBBONS HUNEKER, *Essays* (1929); EDMUND D. JONES (ed.), *English Critical Essays of the Nineteenth Century* (1922); PHYLLIS M. JONES (ed.), *English Critical Essays: Twentieth Century* (1933); WILLIAM PATON KER, *Collected Essays*, 2 vol. (1925, reprinted 1968); KARL MARX and ERICH ENGELS, *Sur la littérature et l'art*, ed. and trans. by JEAN EREVILLE (1936); H.L. MENCKEN, *A Menckens Chrestomathy* (1949); PAUL ELMER MORE, *The Demon of the Absolute* (1928) and *Shellburne Essays*, 11 series (1904-21); ERICH NIETZSCHE, *Ecce Homo and The Birth of Tragedy*, trans. by CLIFTON EADIMAN (1927); HORATIO, *Works*, 10 vol. (1910); GEORGY V. PLEKHANOV, *Art and Society*, trans. by ALERD GOLDSTEIN (1936), a Marxist analysis; EDGAR ALLAN POE, *Selections from the Critical Writings of Edgar Allan Poe*, ed. with an introduction by E.C. PRESCOTT (1909); EZRA POUND, *ABC of Reading* (1934) and *Literary Essays* (1954); HERBERT READ, *Reason and Romanticism* (1926); CHARLES AUGUSTIN SAINTE-BEUVE, *Causeries du lundi*, 15 vol. (1852-62; Eng. trans., 8 vol., 1909-11), contains "What Is a Classic?"; *Shelley's Literary and Philosophical Criticism*, ed. by JOHN SHAWCROSS (1909); LEO TOLSTOY, *What Is Art?*, trans. by AYLMAU MAUDE (1932); LEON TROTSKY, *Literature and Revolution*, trans. by ROSE STRUNSKY (1925, reprinted 1957); PAUL VALÉRY, *Littérature* (1929) and *Variété*, trans. by MALCOLM COWLEY (1927); WILLIAM CARLOS WILLIAMS, *In the American Grain* (1925); EDMUND WILSON, *Axel's Castle* (1931) and *The Triple Thinkers*, rev. ed. (1952 and 1963); EMILIO ZOLA, *The Experimental Novel and Other Essays*, trans. by BELLE M. SHERMAN (1893).

Contemporary: CECIL MAURICE BOWRA, *In General and Particular* (1964); STANLEY BURNISHAW (ed.), *The Poem Itself*, rev. ed. (1967); CYRIL CONNOLLY, *The Modern Movement* (1965); PAUL GOODMAN, *The Structure of Literature* (1954); MARSHALL McLuhan, *Understanding Media* (1964); R.E. SHOLES and R.L. KELLOGG, *The Nature of Narrative* (1966); HERBERT READ, *The Nature of Literature* (1956).

The nature of poetry. The most convenient way to get to know poetry is to read poetry. It would be invidious for the writer of a general article on the subject to prejudice the reader by making a selection of poems or poets; in experience, anyhow, one's acquaintance with poetry comes about chiefly by love and accident, supported, when not undermined, by schools, colleges, and libraries. Beyond that, the bibliographical temptation is to put before the reader numerous learned works that are not poetry but about poetry; whatever their usefulness at various stages of study, and it may be great, they must not substitute for the reading of poetry itself. Therefore no such list is attempted.

The beginning reader, however, may well be able to use some help in interpreting, such as a critical or explanatory anthology. CLEANTH BROOKS and ROBERT PENN WARREN, *Understanding Poetry*, 4th ed. (1976), is still probably the best of its kind, as numerous imitations amply attest. See also TZVETAN TODOROV, *Introduction to Poetics* (1981; originally published in French, 1973), a comprehensive introduction to modern poetics.

Prosody. *Greek and Latin prosody*: PAUL MAAS, *Greek Metre*, trans. by HUGH LLOYD-JONES (1962); ULRICH VON WILAMOWITZ-MOLLFENDORE, *Griechische Verskunst* (1921), the definitive work on the subject but difficult for beginners.

Prose rhythm: MORRIS W. CROLL, *Style, Rhetoric, and Rhythm* (1966), contains classic essays on the period styles of prose and on musical scansion of verse; GEORGE SAINTSBURY, *A History of English Prose Rhythm* (1912, reprinted 1965).

History and uses of English prosody: WILLIAM BEARE, *Latin Verse and European Song: A Study in Accent and Rhythm* (1957); ROBERT BRIDGES, *Milton's Prosody*, rev. ed. (1921); HARVEY S. GROSS, *Sound and Form in Modern Poetry: A Study of Prosody from Thomas Hardy to Robert Lowell* (1964); T.S. OMOND, *English Metrists* (1921, reprinted 1968); GEORGE SAINTSBURY, *A History of English Prosody from the Twelfth Century to the Present*, 3 vol. (1906-10); JOHN THOMPSON, *The Founding of English Metre* (1961).

Theories of prosody: SEYMOUR B. CHATMAN, *A Theory of Meter* (1965); OTTO JESPERSEN, "Notes on Metre," in *The Selected Writings of Otto Jespersen* (1962); WILLIAM K. WIMSATT, JR., and MONROE C. BEARDSLEY, "The Concept of Metre: An Exercise in Abstraction," in WILLIAM K. WIMSATT, JR., *Hateful Contraries* (1965); YVOR WINTERS, "The Audible Reading of Poetry," in *The Function of Criticism* (1957).

Non-Western prosody: ROBERT H. BROWER and EARL MINER, *Japanese Court Poetry* (1961); JAMES LEGGE (ed. and trans.), *The Chinese Classics*, vol. 4 (1960).

General works: PAUL EUSSELL, JR., *Poetic Meter and Poetic Form* (1965); HARVEY S. GROSS (ed.), *The Structure of Verse: Modern Essays on Prosody* (1966); JOSEPH MALOE, *A Manual of English Meters* (1970).

Epic. H.M. CHADWICK, *The Heroic Age* (1912), and H.M. and N.K. CHADWICK, *The Growth of Literature*, vol. 1, *The Ancient Literatures of Europe* (1932), are two classic works on European heroic poetics that are still valuable. A more comprehensive and up-to-date general survey is given in C.M. BOWRA, *Heroic Poetry* (1952); and J. DE VRIES, *Heldenlied und Heldensage* (1961; Eng. trans., *Heroic Song and Heroic Legend*, 1963). A.B. LORD, *The Singer of Tales* (1960), was written by an authority on the Balkan oral epic of the *guslar*. On Homer, see G.S. KIRK, *The Songs of Homer* (1962); W. SCHADFWALDT, *Von Homers Welt und Werk*, 4th ed. (1965); C.M. BOWRA, *Homer and His Forerunners* (1955); and R. CARPENTER, *Folk Tale, Fiction and Saga in the Homeric Epics* (1946). E.R. SCHRODER, *Germanische Heldendichtung* (1935), is a basic reference book for the study of the Germanic epic. J.B. PRITCHARD, *The Ancient Near East* (1958), gives summaries and translations of Akkadian and Ugaritic epics. A theory of the epics of the Indo-Europeans is presented by G. DUMEZIL in *Mythe et épopée*, vol. 1-2 (1968-71). For the use of mythical themes in the epics of the Indo-Europeans, see also D. WARD, *The Divine Twins: An Indo-European Myth in Germanic Tradition* (1968). MICHAEL MURRIN, *The Allegorical Epic* (1980), is a survey of the European tradition.

Fable, parable, and allegory. *Fable*: F. CHAMBRÉ, *Fables* (1927), in Greek and French; S.A. HANDEOD, *Fables of Aesop* (1956); B. PARES, *Krylov's Fables* (1926); MARIANNE MOORE, *Fables of La Fontaine* (1954). For commentary on fables, see: P. CLARAC, *La Fontaine, l'homme et l'oeuvre* (1947); B.F. PERRY, *Aesopica* (1952).

Parable: A.M. HUNTER, *The Parables, Then and Now* (1971); FTA LINNEMANN, *Gleichnisse Jesu*, 3rd ed. (1964; Eng. trans., *The Parables of Jesus*, 1966); T.W. MANSON (ed.), *The Sayings of Jesus as Recorded in the Gospels according to St. Matthew and St. Luke* (1949); D.C. ALLEN, *The Legend of Noah: Renaissance Rationalism in Art, Science and Letters* (1963); HEINZ POLITZER, *Franz Kafka: Parable and Paradox* (1962).

Allegory: (General theory and history): D.C. ALLEN, *Mysteriously Meant: The Rediscovery of Pagan Symbolism and Allegorical Interpretation in the Renaissance* (1970); C.H. DODD, *The Authority of the Bible* (1958); A.S. FLETCHER, *Allegory: The Theory of a Symbolic Mode* (1964); R.M. GRANT, *The Letter and the Spirit* (1958); EDWIN HONIG, *Dark Conceit: The Making of Allegory* (1959); C.S. LEWIS, *The Allegory of Love* (1936); JEAN PEPIN, *Mythe et allégorie* (1958); ROSEMOND TUVE, *Allegorical Imagery* (1966); MAUREEN QUILLIGAN, *The Language of Allegory: Defining the Genre* (1979).

(*Pagan and Christian interpretation*): KENNETH BURKE, *The Rhetoric of Religion* (1961); HENRY CHADWICK, *Early Christian Thought and the Classical Tradition* (1966); C.H. DODD, *The Interpretation of the Fourth Gospel* (1968); A.O. LOVEJOY, *The Great Chain of Being* (1936); H. DE LUBAC, *Exégèse médiévale: les quatre sens de l'Écriture* (1959-64); A. MOMIGLIANO (ed.), *The Conflict Between Paganism and Christianity in the Fourth Century* (1963); G. VON RAD, *Theologie des Alten Testaments*, 2nd ed. (1958; Eng. trans., *Old Testament Theology*), 2 vol., 1962-65); RENE ROQUES, *L'Univers dionysien* (1954); B. SMALLEY, *The Study of the Bible in the Middle Ages*, 2nd ed. (1952); H.A. WOLFSON, *The Philosophy of the Church Fathers*, vol. 1, *Fourth, Trinity, Incarnation* (1956); Philo, *Foundations*

of *Religious Philosophy in Judaism, Christianity, and Islam*, 2 vol. (1947).

(*Typology and typological symbolism*): ERICH AUERBACH, "Figura," in *Scenes from the Drama of European Literature: Six Essays* (1959); A.C. CHARITY, *Events and Their Afterlife: The Dialectics of Christian Typology in the Bible and Dante* (1966); JEAN DANIELOU, *Sacramentum futuri: études sur les origines de la typologie biblique* (1950; Eng. trans., *From Shadows to Reality: Studies in the Biblical Typology of the Fathers*, 1960); AUSTIN EARRER, *A Rebirth of Images: The Making of St. John's Apocalypse* (1949); R.P.C. HANSON, *Allegory and Event* (1959); W.G. MADSEN, *From Shadowy Types to Truth: Studies in Milton's Symbolism* (1968).

(*Medieval allegory*): ERICH AUERBACH, *Dante als Dichter der irdischen Welt* (1929; Eng. trans., *Dante: Poet of the Secular World*, 1961); M.W. BLOOMEFIELD, "Symbolism in Medieval Literature," *Modern Philology*, 56:73-81 (1958), and *Piers Plowman As a Fourteenth-Century Apocalypse* (1962); EDGAR DE BRUYNE, *Études d'esthétique médiévale*, 3 vol. (1946); M.D. CHENU, *La Théologie au douzième siècle* (1957; Eng. trans. of nine selected essays, *Nature, Man, and Society in the Twelfth Century*, 1968); E.R. CURTIUS, *Europäische Literatur und lateinisches Mittelalter* (1948; Eng. trans., *European Literature and the Latin Middle Ages*, 1953); RAYMOND KLIBANSKY, *The Continuity of the Platonic Tradition During the Middle Ages* (1939); C.S. LEWIS, *The Discarded Image: An Introduction to Medieval and Renaissance Literature* (1964); JOSEPH A. MAZZEO, *Medieval Cultural Tradition in Dante's Comedy* (1960); D.W. ROBERTSON and B.F. HUPPE, *Piers Plowman and Scriptural Tradition* (1951); CHARLES SINGLETON, *Dante Studies*, vol. 1, *Commedia* (1954).

(*Renaissance and modern allegory*): DOUGLAS BUSH, *Mythology and the Renaissance Tradition in English Poetry*, rev. ed. (1963); WALTER BENJAMIN, *Ursprung des deutschen Trauerspiels* (1928); HAROLD BLOOM, *The Visionary Company* (1961); A.S. ELETCHER, *The Prophetic Moment: An Essay on Spenser* (1971); ALASTAIR FOWLER, *Triumphal Forms: Structural Patterns in Elizabethan Poetry* (1970); NORTHRUP ERYE, *Fearful Symmetry: A Study of William Blake* (1947); U.M. KAUEMAN, *The Pilgrim's Progress and Traditions in Puritan Meditation* (1966); MICHAEL MURRIN, *The Veil of Allegory: Some Notes Toward a Theory of Allegorical Rhetoric in the English Renaissance* (1969); JEAN SEZNEC, *La Survivance des dieux antiques* (1939; Eng. trans., *The Survival of the Pagan Gods*, rev. ed., 1953); E.M.W. TILLYARD, *The Elizabethan World Picture* (1943); EDGAR WIND, *Pagan Mysteries in the Renaissance*, new ed. (1968).

Ballad. F.J. CHILD (ed.), *The English and Scottish Popular Ballads*, 5 vol. (1882-98), is the canon of traditional balladry; the tunes for which are supplied in B.H. BRONSON (ed.), *Traditional Tunes of the Child Ballads*, 4 vol. (1959-72). Important broadside collections include *The Roxburghe Ballads*, ed. by W. CHAPPELL and J.W. EBSWORTH, 9 vol. (1871-99); and *The Pepys Ballads*, ed. by H.E. ROLLINS, 8 vol. (1929-32). See also *The Common Muse: An Anthology of Popular British Ballad Poetry, XVth-XXth Century*, ed. by V. DE SOLA PINTO and A.E. RODWAY (1957); C.M. SIMPSON, *The British Broadside Ballad and Its Music* (1966); T.P. COFFIN, *The British Traditional Ballad in North America*, rev. ed. (1963); and G.M. LAWS, *Native American Balladry*, rev. ed. (1964). Ballad criticism and scholarship are analyzed in S.B. HUSTVEDT, *Ballad Books and Ballad Men* (1930); D.K. WILGUS, *Anglo-American Folksong Scholarship Since 1898* (1959); A.B. FRIEDMAN, *The Ballad Revival: Studies in the Influence of Popular on Sophisticated Poetry* (1961); C.J. SHARP, *English Folk-Song: Some Conclusions* (1907); G.H. GEROULD, *Ballad of Tradition* (1932); and M.J.C. HODGART, *Ballads* (1950). A.T. QUILLERCOUCH (ed.), *The Oxford Book of Ballads* (1910, reissued 1951); M. LEACH (ed.), *The Ballad Book* (1955); and A.B. FRIEDMAN (ed.), *Folk Ballads of the English Speaking World* (1956), are the standard anthologies.

Romance. Among older works the most notable are RICHARD HURD, *Letters on Chivalry and Romance* (1764); GEORGE ELLIS, *Specimens of Early English Metrical Romances*, 3 vol. (1805); and SIR WALTER SCOTT, "Essay on Romance" in the Supplement to the 1815-24 edition of the *Encyclopædia Britannica*. The academic study of romance as a form of imaginative narrative may be said to have begun in 1897 with the publication of W.P. KER, *Epic and Romance* (2nd ed. 1908, reprinted 1957), and of GEORGE SAINTSBURY, *The Flourishing of Romance and the Rise of Allegory. (Origins and sources)*; EDMOND EARL, *Recherches sur les sources latines des contes et romans courtois du moyen âge* (1913); JESSIE L. WESTON, *From Ritual to Romance* (1920, reprinted 1957); ROGER S. LOOMIS, *Arthurian Tradition and Chrétien de Troyes* (1949); and JEAN MARX, *La Légende arthurienne et le Graal* (1952). (*Nature and development*): FANNI BOGDANOW, *The Romance of the Grail* (1966); EUGENE VINAVER, *The Rise of Romance* (1971); and ROSEMOND TUVE, *Allegorical Imagery* (1966). J.D. BRUCE, *The Evolution of Arthurian Romance*, 2nd ed., 2 vol. (1928), at one time the

standard work in this field, has now been largely superseded by R.S. LOOMIS (ed.), *Arthurian Literature in the Middle Ages: A Collaborative History* (1959). Since 1949 the International Arthurian Society has been publishing an annual *Bibliographical Bulletin* covering the whole range of Arthurian literature in all languages.

Saga. The best guide to current research is the annual *Bibliography of Old Norse-Icelandic Studies* (from 1964); for earlier works on the sagas, see *Islandica* (from 1908). Standard editions of important texts include *Íslenzk fornrit* (from 1933); *Altnordische Saga-Bibliothek*, ed. by G. CEDERSCHILD *et al.*, 18 vol. (1892-1929); *Editiones Arnarnaganae* (from 1958); *Fornaldar Sögur Nordurlanda*, 4 vol. (1950); *Sturlunga saga*, 2 vol. (1946); and *Nelson's Icelandic Texts* (from 1957), with English translations. Useful general surveys of the sagas are PETER HALLBERG, *Den Isländska Sagan* (1956; Eng. trans., *The Icelandic Saga*, 1962); S. NORDAL, *Sagalitteraturen* (1953); and KURT SCHIER, *Sagaliteratur* (1969). For criticism and interpretation of the saga, see WALTER BAETKE, *Über die Entstehung der Isländersagas* (1956); THEODORE M. ANDERSSON, *The Problem of Icelandic Saga Origins* (1964); GABRIEL TURVILLE-PETRE, *Origins of Icelandic Literature* (1953); THEODORE M. ANDERSSON, *The Icelandic Family Saga: An Analytic Reading* (1967); HERMANN PALSSON, *Art and Ethics in Hrafnkel's Saga* (1971); EINAR O. SVEINSSON, *Á Njálsbúð, bok um mikid listaverk* (1943; Eng. trans., *Njals Saga: A Literary Masterpiece*, 1971); GABRIEL TURVILLE-PETRE, *The Heroic Age of Scandinavia* (1951) and *Myth and Religion of the North* (1964); and HERMANN PALSSON and PAUL EDWARDS, *Legendary Fiction in Medieval Iceland* (1970). See also MARGARET SCHLAUCH, *Romance in Iceland* (1934); and E.E. HALVORSEN, *The Norse Version of the Chanson de Roland* (1959).

Of translations into English, the following may be mentioned: L.M. HOLLANDER (ed. and trans.), *Heimskringla: History of the Kings of Norway* (1964) and *The Sagas of Kormák and the Sworn Brothers* (1949); GWYN JONES (ed. and trans.), *The Vatnsdalers' Saga* (1944), *Egil's Saga* (1960), and *Eirik the Red, and Other Icelandic Sagas* (1961); GEORGE JOHNSTON (ed. and trans.), *The Saga of Gisle* (1963); HERMANN PALSSON (ed. and trans.), *Hrafnkel's Saga and Other Icelandic Stories* (1971); PAUL EDWARDS and HERMANN PALSSON (eds. and trans.), *Arrow-Odd: A Medieval Novel* (1970), *Gautrek's Saga, and Other Medieval Tales* (1968), *Hrolf Gautreksson: A Viking Romance* (1972), and *Eyrbyggja Saga* (1973); MAGNUS MAGNUSSON and HERMANN PALSSON (eds. and trans.), *Njal's Saga* (1960), *The Vinland Sagas* (1965), *King Harald's Saga* (1966), and *Laxdaela Saga* (1969); M.H. SCARGILL and MARGARET SCHLAUCH (eds. and trans.), *Three Icelandic Sagas* (1950); J.I. YOUNG (ed. and trans.), *The Prose Edda of Snorri Sturluson: Tales from Norse Mythology* (1954); J.H. MCGREW (ed. and trans.), *Sturlunga Saga*, vol. 1 (1970); DENTON EOX and HERMANN PALSSON (eds. and trans.), *Grettir's Saga* (1973).

Novel. The following works deal in general terms with the reader's approach to the novel: WALTER ALLEN, *Reading a Novel*, rev. ed. (1963); VAN METER AMES, *Aesthetics of the Novel* (1928, reprinted 1966); CLEANTH BROOKS and R.P. WARREN (eds.), *Understanding Fiction*, 3rd. ed. (1979); ALEXANDER COMEORT, *The Novel and Our Time* (1948); PELHAM EDGAR, *The Art of the Novel* (1933, reprinted 1966); WILSON EOLLETT, *The Modern Novel: A Study of the Purpose and Meaning of Fiction*, rev. ed. (1923); E.M. FORSTER, *Aspects of the Novel* (1927, many reprintings); PERCY LUBBOCK, *The Craft of Fiction*, new ed. (1957).

The following are concerned with the problems of writing fiction and are all the work of novelists: PHYLLIS BENTLEY, *Some Observations on the Art of Narrative* (1946); *Conrad's Prefaces to His Works*, with an essay by EDWARD GARNETT (1937); HENRY JAMES, *The Art of Fiction and Other Essays*, ed. by MORRIS ROBERTS (1948), and *The Art of the Novel*, introduction by R.P. BLACKMUR (1934); EDITH WHARTON, *The Writing of Fiction* (1925); THOMAS WOLFE, *The Story of a Novel* (1936).

The various elements of the novel are dealt with in the following: BONAMY DOBREE, *Modern Prose Style*, 2nd ed. (1964); MAREN ELWOOD, *Characters Make Your Story* (1942); MANUEL KOMROEE, *How to Write a Novel* (1950); W. VAN O'CONNOR (ed.), *Forms of Modern Fiction* (1948); GEORGE G. WILLIAMS (ed.), *Readings for Creative Writers* (1938).

The following studies deal with the style and philosophy of the novel in the wider sense: DAVID DAICHES, *The Novel and the Modern World*, rev. ed. (1960); AGNES HANSEN, *Twentieth Century Forces in European Fiction* (1934); ALERED KAZIN, *On Native Grounds* (1942); Y. KRICKORIAN (ed.), *Naturalism and the Human Spirit* (1944); GEORGE LUKACS, *Studies in European Realism* (1950), and *The Historical Novel* (1962); H.J. MULLER, *Modern Fiction* (1937); and MAS'UD ZAVARZADEH, *The Mythopoetic Reality: The Postwar American Nonfiction Novel* (1976).

Short story. K.P. KEMPTON examines the genre, emphasizing

theme and meaning, in *The Short Story* (1947). An excellent analysis of story techniques is offered in SEAN O'FAOLAIN, *The Short Story* (1948). BRANDER MATTHEWS, *The Philosophy of the Short Story* (1901, reprinted 1931), is predicated on Poe's theories. A provocative thesis regarding the nature of stories is presented in FRANK O'CONNOR, *The Lonely Voice: A Study of the Short Story* (1963). H.S. SUMMERS has collected some of the more important discussions of the form in *Discussions of the Short Story* (1963). *Storytellers and Their Art*, ed. by G. SAMPSON and C. BURKHART (1963), contains comments of various authors on the form. More specialized discussions of the form are contained in F.L. PATTEE, *The Development of the American Short Story* (1923); R.B. WEST, *Short Story in America, 1900-1950* (1952); E.K. BENNETT, *A History of the German Novelle*, 2nd ed. rev. by H.M. WAIDSON (1961); and S. TRENKNER, *Greek Novella in the Classical Period* (1958).

Dramatic literature. ALLARDYCE NICOLL, *World Drama* (1949), offers the best survey of the whole field, but should be supplemented by JOHN GASSNER and EDWARD QUINN, *The Reader's Encyclopaedia of World Drama* (1969); and PHYLLIS HARTNOLL, *The Oxford Companion to the Theatre*, 3rd ed. (1967). The classical texts of dramatic theory and criticism may be found in a collection by B.H. CLARK, *European Theories of the Drama*, rev. ed. (1965), which also contains an extensive bibliography. The *Naiya Shastra* of BHARATA, the classic source for Indian dramatic theory, was translated by M.M. GHOSE in 1951.

Books of importance in the development of modern theory on drama are BERNARD BECKERMAN, *Dynamics of Drama* (1970); E.R. BENTLEY, *The Life of the Drama* (1964); KENNETH BURKE, *A Grammar of Motives* (1945); FRANCIS FERGUSON, *The Idea of a Theatre* (1949); S.K. LANGER, *Feeling and Form* (1953); ELDER OLSON, *Tragedy and the Theory of Drama* (1961); RONALD PEACOCK, *The Art of Drama* (1957); J.L. STYAN, *The Elements of Drama* (1960); and KEIR ELAM, *The Semiotics of Theatre and Drama* (1980), a technical semiotic approach.

The finest study of the classical drama of Greece is probably H.D.F. KITTO, *Greek Tragedy*, 3rd ed. (1961); and for the medieval drama are recommended KARL YOUNG, *The Drama of the Medieval Church*, 2 vol. (1933); HARDIN CRAIG, *English Religious Drama of the Middle Ages* (1955); and O.B. HARDISON, *Christian Rite and Christian Drama in the Middle Ages* (1965). Oriental theatre is surveyed in FAUBION BOWERS, *Japanese Theatre* (1952); F.A. LOMBARD, *An Outline History of the Japanese Drama* (1928), which should be read in conjunction with ARTHUR WALEY's classic *The Noh Plays of Japan* (1922); A.C. SCOTT, *The Classical Theatre of China* (1957); A.B. KEITH, *The Sanskrit Drama* (1924); with H.W. WELLS's comparative studies, *The Classical Drama of India* (1963), and *The Classical Drama of the Orient* (1965).

M.C. BRADBROOK, *Themes and Conventions of Elizabethan Tragedy* (1935); and U.M. ELLIS-FERMOR, *Jacobean Drama* (1936), are standard surveys of the English Renaissance drama; and for standard Shakespearean criticism the reader should consult A.M. EASTMAN, *A Short History of Shakespearean Criticism* (1968). The classic source books for the commedia dell'arte are P.L. DUCHARTRE, *La Comédie italienne* (Eng. trans., *The Italian Comedy*, 1929, reprinted 1966); and ALLARDYCE NICOLL, *Masks, Mimes and Miracles* (1931). On the French classical drama H.C. LANCASTER, *A History of French Dramatic Literature in the Seventeenth Century*, 9 vol. (1929-42), is standard; but MARTIN TURNELL, *The Classical Moment* (1947), deals more briefly with Corneille, Racine, and Molière. On Restoration comedy J.L. PALMER, *The Comedy of Manners* (1913); and BONAMY DOBREE, *Restoration Comedy* (1924), remain the best.

American drama is surveyed briefly in W.J. MESERVE, *An Outline History of American Drama* (1965); and A.S. DOWNER, *Fifty Years of American Drama, 1900-1950* (1951). U.M. ELLIS-FERMOR, *The Irish Dramatic Movement*, 2nd ed. (1954), is a comprehensive study of the early years at Dublin's Abbey Theatre; and on Western drama after Ibsen the reader should begin by consulting ERIC BENTLEY, *The Playwright as Thinker* (1946, reprinted 1955); ROBERT BRUSTEIN, *The Theatre of Revolt* (1964); and J.L. STYAN, *The Dark Comedy*, 2nd ed. (1968), an account of the blending of tragic and comic elements in the post-Ibsen theatre.

Comedy. C.L. BARBER, *Shakespeare's Festive Comedy: A Study of Dramatic Form and Its Relation to Social Custom* (1959), Shakespearean comedy considered in relation to archetypal patterns of folk ritual and games; LANF COOPER, *An Aristotelian Theory of Comedy, with an Adaptation of the Poetics and a Translation of the Tractatus Coislinianus* (1922), the only modern text of the *Tractatus*, with an introductory essay relating it to Aristotle's theory of tragedy, and a conjectural reconstruction of the lost treatise on comedy based on the example of the *Poetics*; F.M. CORNFORD, *The Origin of Attic Comedy* (1914; ed. by T.H. GASTER, 1961), an account of the development of Greek comedy from primitive fertility

rites, and of the survival of traces of these ceremonials in the extant plays of Aristophanes; CYRUS HOY, *The Hyacinth Room: An Investigation into the Nature of Comedy, Tragedy, and Tragicomedy* (1964), an examination of the plays of Euripides, Shakespeare, Jonson, Molière, Ibsen, Strindberg, Pirandello, Beckett, and Ionesco; J.W. KRUTCH, *Comedy and Conscience After the Restoration* (1924, reprinted with a new preface and additional bibliographic material, 1949), a study of the decline of Restoration comedy and the rise of sentimental comedy at the end of the 17th and the beginning of the 18th century; K.M. LYNCH, *The Social Mode of Restoration Comedy* (1926), the best available account of the relation of the plays of Dryden, Etherege, Wycherley, Congreve, and their contemporaries to their social milieu; A.W. PICKARD-CAMBRIDGE, *Dithyramb, Tragedy, and Comedy* (1927; 2nd ed. rev. by T.B.L. WEBSTER, 1962), and *The Dramatic Festivals of Athens* (1953; 2nd ed. rev. by J. GOULD and D.M. LEWIS, 1968), the definitive accounts of the origins of Greek comedy and tragedy; and F.H. RISTINE, *English Tragicomedy: Its Origin and History* (1910), the only full-scale account of the subject through the 17th century.

Tragedy. A lengthier development of many of the points made in this section may be found in R.B. SEWALL, *The Vision of Tragedy* (1959). J. JONES, *On Aristotle and Greek Tragedy* (1962), examines the origins of Greek tragedy. Works concentrating on modern tragedy include GEORGE STEINER, *The Death of Tragedy* (1961); WALTER KERR, *Tragedy and Comedy* (1968); and RAYMOND WILLIAMS, *Modern Tragedy* (1966). Special aspects of tragedy are treated in J.M.R. MARGESON, *The Origins of English Tragedy* (1967); EUGENE VINAVER, *Racine and Poetic Tragedy* (1955); and A.C. BRADLEY, *Shakespearean Tragedy* (1904). A useful anthology of writings on tragedy is LIONEL ABEL (ed.), *Moderns on Tragedy* (1967). Other recent works on the subject include RICHMOND HATHORN, *Tragedy, Myth, and Mystery* (1962); MURRAY KRIEGER, *The Tragic Vision: Variations on a Theme in Literary Interpretation* (1960); DOROTHY KROOK, *Elements of Tragedy* (1969); and TIMOTHY J. REISS, *Tragedy and Truth* (1980).

Satire. DAVID WORCESTER, *The Art of Satire* (1940), a study of rhetorical techniques available to the satirist; JAMES R. SUTHERLAND, *English Satire* (1958), a sound scholarly history; ALVIN B. KERNAN, *The Cankered Muse: Satire of the English Renaissance* (1959), valuable theory and criticism; ROBERT C. ELLIOTT, *The Power of Satire: Magic, Ritual, Art* (1960), on the origins of satire in magic and its development into an art; RONALD PAULSON, *The Fictions of Satire* (1967), a study of satire in fiction from Lucian to Swift, and (ed.), *Satire: Modern Essays in Criticism* (1971), an authoritative and indispensable collection; MATTHEW J.C. HODGART, *Satire* (1969), a well-illustrated, readable survey of satire in many forms and in many countries.

Other genres. On the essay, see R.D. O'LEARY, *The Essay* (1928), which analyzes the essay theoretically and examines several categories of essayists. DAVID DAICHES gives a brief and lively introduction to his anthologies of essays: *A Century of the Essay: British and American* (1951), and *More Literary Essays* (1968). On letter writing in general, the best piece is by GUSTAVE LANSON, in French: Introduction to *Choix de lettres du XVII^e siècle*, 5th rev. ed. (1898), reprinted in Lanson's *Essais de méthode, de critique et d'histoire littéraire*, pp. 243-258 (1965). On personal literature, the diary, autobiography, and the questions raised by those forms of prose, see HENRI PEYRE, *Literature and Sincerity* (1963).

Biography. *Critical and scholarly books:* JAMES L. CLIFFORD, *From Puzzles to Portraits: Problems of a Literary Biographer* (1970), examples of the author's own research followed by an analysis of biographical problems; LEON EDEL, *Literary Biography* (1959), essentially an account of the methods, psychological and narrative, used by the author in his multivolume life of Henry James; JOHN A. GARRATY, *The Nature of Biography* (1957), a historical survey coupled with a study of biographical methods, with emphasis on aids offered by psychology; PAUL M. KENDALL, *The Art of Biography* (1965), a historical survey, with emphasis on contemporary biography, and a study of biographical problems from the viewpoint of a practicing biographer; ANDRE MAUROIS, *Aspects de la biographie* (1928; Eng. trans. 1930) and HAROLD NICOLSON, *The Development of English Biography* (1928), particularly interesting for complementary views of the "new" biography of the 1920s by two eminent biographers; ROY PASCAL, *Design and Truth in Autobiography* (1960), a historical survey and a study of the chief problems, aspects, and varieties of autobiography; WILLIAM M. RUNYAN, *Life Histories and Psychobiography* (1982), a discussion of methodologies used in conducting psychobiographical research.

Anthologies: JAMES L. CLIFFORD (ed.), *Biography as an Art: Selected Criticism 1560-1960* (1962); WILLIAM H. DAVENPORT and BEN SIEGEL (eds.), *Biography Past and Present* (1965), contains a number of critical essays as well as biographical

selections; EDGAR JOHNSON (ed.), *A Treasury of Biography* (1941); JOHN C. METCALFE (ed.), *The Stream of English Biography* (1930).

Literary criticism. A useful compilation of essential texts is MARK SCHORER, JOSEPHINE MILES, and GORDON MCKENZIE (eds.), *Criticism*, rev. ed. (1958). The best survey of critical history is WILLIAM K. WIMSATT, JR., and CLEANTH BROOKS, *Literary Criticism: A Short History* (1957); G.M.A. GRUBE, *The Greek and Roman Critics* (1965); JOEL E. SPINGARN, *A History of Literary Criticism in the Renaissance*, 5th ed. (1925, paperback edition 1963); WALTER J. BATE, *From Classic to Romantic* (1946, reprinted 1961); and RENE WELLEK, *A History of Modern Criticism, 1750-1950*, 4 vol. (1955-65), are more specialized historical studies. Important theoretical statements are M.H. ABRAMS, *The Mirror and the Lamp* (1953); RENE WELLEK and AUSTIN WARREN, *Theory of Literature*, 3rd rev. ed. (1966); NORTHROP FRYE, *Anatomy of Criticism* (1957); and WAYNE C. BOOTH, *The Rhetoric of Fiction* (1961). WILLIAM EMPSON, *Seven Types of Ambiguity*, 3rd ed. (1956, reprinted 1963); ERICH AUERBACH, *Mimesis* (1946; Eng. trans. 1953); and LIONEL TRILLING, *The Liberal Imagination* (1950), are representative examples of modern criticism, combining theory with analysis of a wide variety of texts. See also DOUWE W. FOKKEMA and ELRUD KUNNE-IBSCH, *Theories of Literature in the Twentieth Century: Structuralism, Marxism, Aesthetics of Reception, Semiotics* (1978).

Children's literature. *Historical, critical:* (Europe): BETTINA HURLIMANN, *Europäische Kinderbücher in drei Jahrhunderten*, 2nd ed. (1963; Eng. trans., *Three Centuries of Children's Books in Europe*, 1968). (England): GILLIAN AVERY, *Nineteenth Century Children: Heroes and Heroines in English Children's Stories 1780-1900* (1965); FLORENCE V. BARRY, *A Century of Children's Books* (1922, reprinted 1968); MARCUS CROUCH, *Treasure Seekers and Borrowers: Children's Books in Britain 1900-1960* (1962); F.J. HARVEY DARTON, *Children's Books in England: Five Centuries of Social Life* (1932); ROGER LANCELYN GREEN, *Tellers of Tales: British Authors of Children's Books from 1800 to 1964*, rev. ed. (1965); PERCY MUIR, *English Children's Books, 1600 to 1900* (1954); M.F. THWAITE, *From Primer to Pleasure: An Introduction to the History of Children's Books in England, from the Invention of Printing to 1900* (1963); JOHN ROWE TOWNSEND, *Written for Children: An Outline of English Children's Literature* (1965). (Canada): SHEILA EGGOFF, *The Republic of Childhood: A Critical Guide to Canadian Children's Literature in English* (1967). (Anglo-American mainly): CORNELIA MEIGS et al., *A Critical History of Children's Literature: A Survey of Children's Books in English from Earliest Times to*

the Present, rev. ed. (1969). (Germany): IRENE DYHRENFURTH-GRAEBSCH, *Geschichte des deutschen Jugendbuches*, 3rd rev. ed. (1967); H.L. KOSTER, *Geschichte der deutschen Jugendliteratur* (1968). (Sweden): EVA VON ZWEIGBERGK, *Barnboken I Sverige 1750-1950* (1965). (France): MARIE-THERESE LATZARUS, *La Littérature enfantine en France dans la seconde moitié du XIX^e siècle* (1923); JEAN DE TRIGON, *Histoire de la littérature enfantine de ma Mère l'Oye au Roi Babar* (1950). (Italy): PIERO BARGELLINI, *Canto alle rondini: panorama storico della letteratura infantile*, 6th ed. (1967); GIUSEPPE FANCIULLI, *Scrittori per l'infanzia*, 3rd ed. (1968); LOUISE RESTIEAUX HAWKES, *Before and After Pinocchio: A Study of Italian Children's Books* (1933). (Spain): CARMEN BRAVO VILLASANTE (ed.), *Historia de la literatura infantil española* (1963). (Latin America): CARMEN BRAVO VILLASANTE, *Historia y antología de la literatura infantil iberoamericana*, 2 vol. (1966); DORA PASTORIZA DE ETCHEBARNE, *El cuento en la literatura infantil, ensayo crítico* (1962).

General: RICHARD BAMBERGER, *Jugendlektüre*, 2nd ed. (1965); ELEANOR CAMERON, *The Green and Burning Tree: On the Writing and Enjoyment of Children's Books* (1969); KORNEI CHUKOVSKY, *From Two to Five*, rev. ed. (1968; Eng. trans. of the 20th Russian ed. of 1968); HANS CORNIOLEY, *Beiträge zur Jugendbuchkunde* (1966); MARGERY FISHER, *Intent upon Reading: A Critical Appraisal of Modern Fiction for Children* (1961); PAUL HAZARD, *Les Livres, les enfants et les hommes* (1932; Eng. trans., *Books, Children and Men*, 4th ed., 1960); ENZO PETRINI, *Avviamento critico alla letteratura giovanile* (1958); LILLIAN H. SMITH, *The Unreluctant Years: A Critical Approach to Children's Literature* (1953); DOROTHY M. WHITE, *Books Before Five* (1954); KAY E. VANDERGRIFT, *Child and Story* (1981).

Bibliographic: VIRGINIA HAVILAND, *Children's Literature: A Guide to Reference Sources* (1966); ANNE PELLOWSKI, *The World of Children's Literature* (1968).

Biographical: BRIAN DOYLE (ed.), *The Who's Who of Children's Literature* (1968); MURIEL FULLER (ed.), *More Junior Authors* (1963); STANLEY J. KUNITZ and HOWARD HAYCRAFT (eds.), *The Junior Book of Authors*, 2nd ed. rev. (1951).

Illustration: BETTINA HURLIMANN, *Die Welt im Bilderbuch* (1965; Eng. trans., *Picture-Book World*, 1968); LEE KINGMAN, JOANNA FOSTER, and RUTH GILES LONTOFT (comps.), *Illustrators of Children's Books: 1957-1966* (1968); DIANA KLEMIN, *The Art of Art for Children's Books: A Contemporary Survey* (1966); BERTHA E. MAHONY et al. (comps.), *Illustrators of Children's Books 1744-1945* (1947); BERTHA MAHONY MILLER et al. (comps.), *Illustrators of Children's Books, 1946-1956* (1958).

The History of Western Literature

Diverse as they are, European literatures, like European languages, are parts of a common heritage. Greek, Latin, Germanic, Baltic and Slavic, Celtic, and Romance languages are all members of the Indo-European family. (Finnish and Hungarian and Semitic languages of the eastern Mediterranean, such as Hebrew, are not Indo-European. Literatures in these languages are, however, closely associated with major Western literatures and are often included among them.) The common literary heritage is essentially that originating in ancient Greece and Rome. It was preserved, transformed, and spread by Christianity and thus transmitted to the vernacular languages of the European Continent, the Western Hemisphere, and other regions that were settled by Europeans. To the present day, this body of writing displays a unity in its main features that sets it apart from the literatures of the rest of the world. Such common characteristics are considered here.

For specific information about the major national literatures or literary traditions of the West, see such *Macropædia* articles as AMERICAN LITERATURE; ENGLISH LITERATURE; GERMAN LITERATURE; GREEK LITERATURE; LATIN-AMERICAN LITERATURE; and SCANDINAVIAN LITERATURE. Various other Western literatures—including those in the Armenian, Bulgarian, Estonian, Lithuanian, and Romanian languages—are treated in separate entries of the *Micropædia*.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, Part Six, Division II, Section 621.

The article is divided into the following sections:

Ancient literature	216
Medieval literature	216
The Renaissance	217
The 17th century	217
The 18th century	218
The 19th century	218
Romanticism	
Post-Romanticism	
The 20th century	219
Bibliography	220

ANCIENT LITERATURE

The stark fact about ancient Western literature is that the greater part of it has perished. Some of it had been forgotten before it was possible to commit it to writing; fire, war, and the ravages of time have robbed posterity of most of the rest; and the restitutions that archaeologists and paleographers achieve from time to time are small. Yet surviving writings in Greek and far more in Latin have included those that on ancient testimony marked the heights reached by the creative imagination and intellect of the ancient world.

Five ancient civilizations, Babylon and Assyria, Egypt, Greece, Rome, and the culture of the Israelites in Palestine, each came into contact with one or more of the others. The two most ancient, Assyro-Babylonia, with its broken clay tablets, and Egypt, with its rotted papyrus rolls, make no direct literary signal to the modern age; yet Babylon produced the first full code of laws and two epics of archetypal myth, which came to be echoed and re-echoed in distant lands, and Egypt's mystical intuition of a supernatural world caught the imagination of the Greeks and Romans. Hebrew culture exerted its greatest literary influence on the West because of the place held by its early writings as the Old Testament of the Christian Bible; and this literature profoundly influenced Western consciousness through translation from about the time

of St. Augustine onward into every vernacular language as well as into Latin. Until then, Judaism's concentrated spirituality set it apart from the Greek and Roman world.

Though influenced by the religious myths of Mesopotamia, Asia Minor, and Egypt, Greek literature has no direct literary ancestry and appears as self-originated. Roman writers looked to Greek precept for themes, treatment, and choice of verse and metre. Rome eventually passed the torch on to the early Middle Ages, by which time Greek had been subsumed under a wholly Latin tradition and was only rediscovered in its own right at the Renaissance—the "classical" tradition afterward becoming a threat to natural literary development, particularly when certain critics of the 17th century began to insist that the subjects and style of contemporary writing should conform with those employed by Greece and Rome.

All of the chief kinds of literature—epic, tragedy, comedy, lyric, satire, history, biography, and prose narrative—were established by the Greeks and Romans, and later developments have for the most part been secondary extensions. The Greek epic of Homer was the model for the Latin of Virgil; the lyric fragments of Alcaeus and Sappho were echoed in the work of Catullus and Ovid; the history of Thucydides was succeeded by that of Livy and Tacitus; but the tragedy of the great Athenians of the 5th century BC had no worthy counterpart in Roman Seneca nor had the philosophical writings of Plato and Aristotle in those of any ancient Roman, for the practical Romans were not philosophers. Whereas Greek writers excelled in abstraction, the Romans had an unusually concrete vision and, as their art of portraiture shows, were intensely interested in human individuality.

In sum, the work of these writers and others and perhaps especially that of Greek authors expresses the imaginative and moral temper of Western man. It has helped to create his values and to hand on a tradition to distant generations. Homer's epics extend their concern from the right treatment of strangers to behaviour in situations of deep involvement among hero rivals, their foes, and the over-seeing gods; the tragedies of Aeschylus and Sophocles are a sublime expression of man's breakthrough into moral awareness of his situation. Among Roman authors an elevated Stoicism stressing the sense of duty is common to many, from Naevius, Ennius, and Cato to Virgil, Horace, and Seneca. A human ideal is to be seen in the savage satire of Juvenal and in Anacreon's songs of love and wine, as it is in the philosophical thought of Plato and Aristotle. It is given voice by a chorus of Sophocles, "Wonders are many, but none is more wonderful than man, the power that crosses the white sea. . . ." The human ideal held up in Greek and Latin literature, formed after civilization had emerged from earlier centuries of barbarism, was to be transformed, before the ancient world came to its close, into the spiritual ideal of Judeo-Christianity, whose writers foreshadowed medieval literature.

MEDIEVAL LITERATURE

Medieval, "belonging to the Middle Ages," is used here to refer to the literature of Europe and the eastern Mediterranean from as early as the establishment of the Eastern Roman, or Byzantine, Empire about AD 300 for medieval Greek, from the period following upon the fall of Rome in 476 for medieval Latin, and from about the time of Charlemagne and the Carolingian Renaissance he fostered in France (c. 800) to the end of the 15th century for most written vernacular literatures.

The establishment of Christianity throughout the territories that had formed the Roman Empire meant that Europe was exposed to and tutored in the systematic approach to life, literature, and religion developed by the early Church Fathers. In the West, the fusion of Christian

and classical philosophy formed the basis of the medieval habit of interpreting life symbolically. Through St. Augustine, Platonic and Christian thought were reconciled: the permanent and uniform order of the Greek universe was given Christian form; nature became sacramental, a symbolic revelation of spiritual truth. Classical literature was invested with this same symbolism; exegetical, or interpretative, methods first applied to the Scriptures were extended as a general principle to classical and secular writings. The allegorical or symbolic approach that found in Virgil a pre-Christian prophet and in the *Aeneid* a narrative of the soul's journey through life to paradise (Rome) belonged to the same tradition as Dante's allegorical conception of himself and his journey in *The Divine Comedy*.

The role of the church in preserving ancient literatures

The church not only established the purpose of literature but preserved it. St. Benedict's monastery at Monte Cassino in Italy was established in 529, and other monastic centres of scholarship followed, particularly after the 6th- and 7th-century Irish missions to the Rhine and Great Britain and the Gothic missions up the Danube. These monasteries were able to preserve the only classical literature available in the West through times when Europe was being raided by Goths, Vandals, Franks, and, later, Norsemen in succession. The classical Latin authors so preserved and the Latin works that continued to be written predominated over vernacular works throughout most of the period. St. Augustine's *City of God*, the Venerable Bede's *Ecclesiastical History*, the Danish chronicle of Saxo Grammaticus, for example, were all written in Latin, as were most major works in the fields of philosophy, theology, history, and science.

The main literary values of the period are found in vernacular works. The pre-Christian literature of Europe belonged to an oral tradition that was reflected in the *Poetic Edda* and the sagas, or heroic epics, of Iceland, the Anglo-Saxon *Beowulf*, and the German *Song of Hildebrand*. These belonged to a common Germanic alliterative tradition, but all were first recorded by Christian scribes at dates later than the historical events they relate, and the pagan elements they contain were fused with Christian thought and feeling. The mythology of Icelandic literature was echoed in every Germanic language and clearly stemmed from a common European source. Only the Scandinavian texts, however, give a coherent account of the stories and personalities involved. Numerous ballads in different countries also reflect an earlier native tradition of oral recitation. Among the best known of the many genres that arose in medieval vernacular literatures were the romance and the courtly love lyric, both of which combined elements from popular oral traditions with those of more scholarly or refined literature and both derived largely from France. The romance used classical or Arthurian sources in a poetic narrative that replaced the heroic epics of feudal society, such as *The Song of Roland*, with a chivalrous tale of knightly valour. In the romance, complex themes of love, loyalty, and personal integrity were united with a quest for spiritual truth, an amalgam that was represented in every major western European literature of the time. The love lyric has had a similarly heterogeneous background. The precise origins of courtly love are disputed, as is the influence of a popular love poetry tradition; it is clear, however, that the idealized lady and languishing suitor of the poets of southern and northern France were imitated or reinterpreted throughout Europe—in the Sicilian school of Italy, the minnesingers (love poets) of Germany, and in a Latin verse collection, *Carmina Burana*.

Medieval drama began in the religious ceremonies that took place in church on important dates in the Christian calendar. The dramatic quality of the religious service lent itself to elaboration that perhaps first took the form of gestures and mime and later developed into dramatic interpolations on events or figures in the religious service. This elaboration increased until drama became a secular affair performed on stages or carts in town streets or open spaces. The players were guild craftsmen or professional actors and were hired by towns to perform at local or religious festivals. Three types of play developed: the mystery, the miracle, and the morality. The titles and themes

of medieval drama remained religious but their pieces' titles can belie their humorous or farcical and sometimes bawdy nature. One of the best known morality plays was translated from Dutch to be known in English as *Everyman*. A large majority of medieval literature was anonymous and not easily dated. Some of the greatest figures—Dante, Chaucer, Petrarch, and Boccaccio—came late in the period, and their work convincingly demonstrates the transitional nature of the best of medieval literature, for, in being master commentators of the medieval scene, they simultaneously announced the great themes and forms of Renaissance literature.

THE RENAISSANCE

The name Renaissance ("Rebirth") is given to the historical period in Europe that succeeded the Middle Ages. The awakening of a new spirit of intellectual and artistic inquiry, which was the dominant feature of this political, religious, and philosophical phenomenon, was essentially a revival of the spirit of ancient Greece and Rome; in literature this meant a new interest in and analysis of the great classical writers. Scholars searched for and translated "lost" ancient texts, whose dissemination was much helped by developments in printing in Europe from about 1450.

Art and literature in the Renaissance reached a level unattained in any previous period. The age was marked by three principal characteristics: first, the new interest in learning, mirrored by the classical scholars known as humanists and instrumental in providing suitable classical models for the new writers; second, the new form of Christianity, initiated by the Protestant Reformation led by Martin Luther, which drew men's attention to the individual and his inner experiences and stimulated a response in Catholic countries summarized by the term Counter-Reformation; third, the voyages of the great explorers that culminated in Christopher Columbus' discovery of America in 1492 and that had far-reaching consequences on the countries that developed overseas empires, as well as on the imaginations and consciences of the most gifted writers of the day.

To these may be added many other factors, such as the developments in science and astronomy and the political condition of Italy in the late 15th century. The new freedom and spirit of inquiry in the Italian city-states had been a factor in encouraging the great precursors of the Renaissance in Italy, Dante, Petrarch, and Boccaccio. The flowering of the Renaissance in France appeared both in the poetry of the poets making up the group known as the *Pléiade* and in the reflective essays of Michel de Montaigne, while Spain at this time produced its greatest novelist, Miguel de Cervantes. Another figure who stood out above his contemporaries was the Portuguese epic poet Luís Camões, while drama flourished in both Spain and Portugal, being represented at its best by Lope de Vega and Gil Vicente. In England, too, drama dominated the age, a blend of Renaissance learning and native tradition lending extraordinary vitality to works of Christopher Marlowe, Ben Jonson, John Webster, and others, while Shakespeare, England's greatest dramatic and poetic talent, massively spanned the end of the 16th century and the beginning of the 17th.

In the 16th century the Dutch scholar Desiderius Erasmus typified the development of humanism, which embodied the spirit of critical inquiry, regard for classical learning, intolerance of superstition, and high respect for man as God's most intricate creation. An aspect of the influence of the Protestant Reformation on literature was the number of great translations of the Bible, including an early one by Erasmus, into vernacular languages during this period, setting new standards for prose writing. The impetus of the Renaissance carried well into the 17th century, when John Milton reflected the spirit of Christian humanism.

THE 17TH CENTURY

The 17th century was a period of unceasing disturbance and violent storms, no less in literature than in politics and society. The Renaissance had prepared a receptive environment essential to the dissemination of the ideas of the new science and philosophy. The great question of the

century, which confronted serious writers from Donne to Dryden, was Michel de Montaigne's "What do I know?" or, in expanded terms, the ascertainment of the grounds and relations of knowledge, faith, reason, and authority in religion, metaphysics, ethics, politics, economics, and natural science.

The questioning attitude that characterized the period is seen in the works of its great scientists and philosophers: Descartes's *Discourse on Method* (1637) and Pascal's *Pensées* (written 1657–58) in France; Bacon's *Advancement of Learning* (1605) and Hobbes's *Leviathan* (1651) in England. The importance of these works has lain in their application of a skeptical, rationalist mode of thought not only to scientific problems but to political and theological controversy and general problems of understanding and perception. This fundamental challenge to both thought and language had profound repercussions in man's picture of himself and was reflected in what T.S. Eliot described as "the dissociation of sensibility," which Eliot claimed took root in England after the Civil War, whereby, in contrast to the Elizabethan and Jacobean writers who could "devour any kind of experience," later poets in English could not think and feel in a unified way.

A true picture of the period must also take into account the enormous effect of social and political upheavals during the early and middle parts of the century. In England, where the literary history of the period is usually divided into two parts, the break seems to fall naturally with the outbreak of the Civil War (1642–51), marked by a closure of the theatres in 1642, and a new age beginning with the restoration of the monarchy in 1660. In France the bitter internecine struggle of the Fronde (1648–53) similarly divided the century and preceded possibly the greatest period of all French literature—the age of Molière, Racine, Boileau, and La Fontaine. In Germany the early part of the century was dominated by the religious and political conflicts of the Thirty Years' War (1618–48) and thereafter by the attempts of German princes to emulate the central power and splendour of Louis XIV's French court at Versailles. The Netherlands was also involved in the first part of the century in a struggle for independence from Spain (the Eighty Years' War, 1568–1648) that resulted not only in the achievement of this but also in the "Golden Age" of Dutch poetry—that of Henric Spieghel, Daniël Heinsius, and Gerbrand Bredero.

The civil, political, and religious conflicts that dominated the first half of the century were in many ways also the characteristic response of the Counter-Reformation. The pattern of religious conflict was reflected in literary forms and preoccupations. One reaction to this—seen particularly in Italy, Germany, and Spain but also in France and England—was the development of a style in art and literature known as Baroque. This development manifested itself most characteristically in the works of Giambattista Marino in Italy, Luis de Góngora in Spain, and Martin Opitz in Germany. Long regarded by many critics as decadent, Baroque literature is now viewed in a more favourable light and is understood to denote a style the chief characteristics of which are elaboration and ornament, the use of allegory, rhetoric, and daring artifice.

If Baroque literature was the characteristic product of Italy and Germany in this period, Metaphysical poetry was the most outstanding feature in English verse of the first half of the century. This term, first applied by Dryden to John Donne and expanded by Dr. Johnson, is now used to denote a range of poets who varied greatly in their individual styles but who possessed certain affinities with Baroque literature, especially in the case of Richard Crashaw.

Perhaps the most characteristic of all the disputes of the 17th century was that in which the tendency to continue to develop the Renaissance imitation of the classics came into conflict with the aspirations and discoveries of new thinkers in science and philosophy and new experimenters with literary forms. In France this appeared in a struggle between the Ancients and Moderns, between those who thought that literary style and subject should be modeled on classical Greek and Latin literature and supporters of native tradition. In Spain a similar conflict was expressed

in a tendency toward ornament, Latinization, and the classics (*culteranismo*) and that toward a more concise, profound, and epigrammatic style (*conceptismo*). This conflict heralded through the Moderns in France and the idea of *conceptismo* in Spain a style of prose writing suitable to the new age of science and exploration. The Moderns in France were largely, therefore, followers of Descartes. In England a similar tendency was to be found in the work of the Royal Society in encouraging a simple language, a closer, naked, natural way of speaking, suitable for rational discourse, paralleled by the great achievements in prose of John Milton and John Dryden.

THE 18TH CENTURY

To call the 18th century the Age of Reason is to seize on a useful half-truth but to cause confusion in the general picture, because the primacy of reason had also been a mark of certain periods of the previous age. It is more accurate to say that the 18th century was marked by two main impulses: reason and passion. The respect paid to reason was shown in pursuit of order, symmetry, decorum, and scientific knowledge; the cultivation of the feelings stimulated philanthropy, exaltation of personal relationships, religious fervour, and the cult of sentiment, or sensibility. In literature the rational impulse fostered satire, argument, wit, plain prose; the other inspired the psychological novel and the poetry of the sublime.

The cult of wit, satire, and argument is evident in England in the writings of Alexander Pope, Jonathan Swift, and Samuel Johnson, continuing the tradition of Dryden from the 17th century. The novel was established as a major art form in English literature partly by a rational realism shown in the works of Henry Fielding, Daniel Defoe, and Tobias Smollett and partly by the psychological probing of the novels of Samuel Richardson and of Laurence Sterne's *Tristram Shandy*. In France the major characteristic of the period lies in the philosophical and political writings of the Enlightenment, which had a profound influence throughout the rest of Europe and foreshadowed the French Revolution. Voltaire, Jean-Jacques Rousseau, Charles de Montesquieu, and the Encyclopédistes Denis Diderot and Jean d'Alembert all devoted much of their writing to controversies about social and religious matters, often involving direct conflict with the authorities. In the first part of the century, German literature looked to English and French models, although innovative advances were made by the dramatist and critic Gotthold Ephraim Lessing. The great epoch of German literature came at the end of the century, when cultivation of the feelings and of emotional grandeur found its most powerful expression in what came to be called the *Sturm und Drang* ("Storm and Stress") movement. Associated with this were two of the greatest names of German literature, Johann Wolfgang von Goethe and Friedrich Schiller, both of whom in drama and poetry advanced far beyond the turbulence of *Sturm und Drang*.

THE 19TH CENTURY

The 19th century in Western literature—one of the most vital and interesting periods of all—has special interest as the formative era from which many modern literary conditions and tendencies derived. Influences that had their origins or were in development in this period—Romanticism, Symbolism, Realism—are reflected in the current of modern literature, and many social and economic characteristics of the 20th century were determined in the 19th.

Romanticism. The predominant literary movement of the early part of the 19th century was Romanticism, which in literature had its origins in the *Sturm und Drang* period in Germany. An awareness of this first phase of Romanticism is an important correction to the usual idea of Romantic literature as something that began in English poetry with William Wordsworth and Samuel Taylor Coleridge and the publication of *Lyrical Ballads* in 1798. Moreover, although it is true that the French Revolution of 1789 and the Industrial Revolution were two main political and social factors affecting the Romantic poets of early 19th-century England, many characteristics of Romanticism in literature sprang from literary or philosophical sources. A

The effects of war on 17th-century literature

Baroque literature

The conflict between the classics and the new science

philosophical background was provided in the 18th century chiefly by Jean-Jacques Rousseau, whose emphasis on the individual and the power of inspiration influenced Wordsworth and also such first-phase Romantic writers as Friedrich Hölderlin and Ludwig Tieck in Germany and the French writer Bernardin de Saint-Pierre, whose *Paul et Virginie* (1787) anticipated some of the sentimental excesses of 19th-century Romantic literature. Positive as it was, the influence of Rousseau must also be seen as a partly negative reaction against 18th-century rationalism with its emphasis on intellect.

Belief in self-knowledge was, indeed, a principal article of Romantic faith. Late 18th-century French writers such as Fabre d'Olivet sought to explain the physical world by an idea of a "breath of life" similar to the "inspiration" of Wordsworth and Coleridge. The Romantics believed that the real truth of things could be explained only through examination of their own emotions in the context of nature and the primitive. Because of this emphasis on inspiration, the poet came to assume a central role—that of seer and visionary. Simultaneously, such formal conventions as imitation of the classics were rejected as binding rules. A new directness of the poet's role emphasized the language of the heart and of ordinary men, and Wordsworth even tried to invent a new simplified diction. Poetry became divorced from its 18th-century social context, and a poet was answerable only to ultimate truth and himself. Two classic poses of the Romantic poet were the mystic visionary of John Keats and the superman of Lord Byron—indeed, satirization of the Byronic hero was to become a theme of later novelists such as Fyodor Dostoyevsky, even though he himself had Romantic antecedents.

The fact that Dostoyevsky was a Russian showed how the Romantic stream flowed across Europe. In Spain and Italy, Hungary, Poland, and the Balkans, it took the form of drama, which in England failed to produce great works. The early and middle 19th century was a time of poetry and prose rather than of drama. The Romantic style in poetry was seen everywhere in Europe—in José de Espronceda in Spain; Ugo Foscolo and Giacomo Leopardi in Italy, where it became identified with nationalist sentiments; Aleksandr Pushkin and Mikhail Lermontov in Russia; Adam Mickiewicz in Poland. In America, a Romantic thread also allied with the emergence of national feeling could be seen in the adventurous stories of James Fenimore Cooper; in the supernatural and mystic element in Edgar Allan Poe; in the poetry of Walt Whitman and Henry Wadsworth Longfellow; and in the Transcendentalist theories of Ralph Waldo Emerson and Henry David Thoreau, which, as Wordsworth's pronouncements had done, affirmed the power of "insight" to transcend ordinary logic and experience.

The impetus of Romantic poetry began to slacken after about 1830 and gave way to more objective styles, although many of its themes and devices, such as the misunderstood artist or the unhappy lover, continued to be employed.

Post-Romanticism. Arguably the first post-Romantic poet was a German, Heinrich Heine, but German poetry in the mid-19th century mostly followed Wordsworth, though new tendencies were to be found in August von Platen Hallermünde and an Austrian, Nikolaus Lenau. The principal development was to be seen in France in the growth of a movement known as Parnassianism. Originating with Théophile Gautier, Parnassianism in some ways was an offshoot of Romanticism rather than a reaction against it. In concentrating on the purely formal elements of poetry, on aesthetics, and on "art for art's sake," it changed the direction of French poetry and had much influence abroad. Its most illustrious representative was Charles Baudelaire, who believed that "everything that is not art is ugly and useless." Another branch of new development was the growth of Impressionism and the Symbolist movement, a result of "borrowing" from movements in painting, sculpture, and music. Paul Verlaine, foremost of the Impressionists, used suggestion, atmosphere, and fleeting rhythms to achieve his effects. Symbolism, a selective use of words and images to evoke tenuous moods and meanings, is conveyed in the work of

Stéphane Mallarmé and Arthur Rimbaud. The advance of French poetry in the middle and later part of the century was an achievement of individuals, based on invention of a personal idiom.

The spread of education and, in England, of circulating libraries increased a demand for novels. At the beginning of the 19th century Jane Austen had already satirized the excesses of the Gothic novel, a harbinger of medievalizing Romanticism in the latter part of the 18th century, in *Northanger Abbey* and the conflict of sense and Romantic sensibility in *Sense and Sensibility*. In France the conflict of intelligence and emotion appeared in the work of Benjamin Constant (*Adolphe*, 1816) and most notably in *Le Rouge et le noir* (1830) of Stendhal and later in Gustave Flaubert's *Madame Bovary* (1857). The detailed verbal scrupulousness and Realism exhibited in the work of Flaubert and of Honoré de Balzac were carried forward by Guy de Maupassant in France and Giovanni Verga in Italy; they culminated in the extreme Naturalism of Émile Zola, who described his prose in novels such as *Thérèse Raquin* (1867) as "literary surgical autopsy."

But Realism and nationalism seem irrelevant as descriptions of the great writers of the period—for example, George Eliot, Charles Dickens, and Thomas Hardy in England and Nikolay Gogol, Ivan Turgenev, Leo Tolstoy, Fyodor Dostoyevsky, and Anton Chekhov in Russia. In such writers there was a distinct bias toward literature with a social purpose, stimulated by awakening forces of liberalism, humanism, and socialism in many Western countries.

A decline of the Romantic theatre into melodrama was fairly general in Europe, and it was slower than the novel to take up problems of contemporary life. When revival came, through the work of a Norwegian, Henrik Ibsen, Romantic conflicts of visionary and realist, individual and society were restated, and this was true also of the plays of August Strindberg in Sweden. In Russia a modern theatre became a vital influence that could trace its beginnings back to Gogol's *Government Inspector* (1836) but was to be felt later in the century in Turgenev's *Month in the Country* (1850) and, above all, in the work of Anton Chekhov, a great dramatist of the period.

THE 20TH CENTURY

When the 20th century began, social and cultural conditions that prevailed in Europe and America were not too different from those of the middle and late 19th century. Continuity could be seen, for example, in the work of four novelists writing in English at the turn of the century and after. Joseph Conrad, Thomas Hardy, Henry James, and D.H. Lawrence all demonstrated in the progress of their work the transition from a relatively stable world at the end of the 19th century to a new age that began with World War I. The awakening of a new consciousness in literature was also to be traced in such works of fiction as the first volume of Marcel Proust's *Remembrance of Things Past* (*Swann's Way*, 1913), André Gide's *Vatican Cellars* (1914), James Joyce's *Ulysses* (1922), Franz Kafka's *Trial* (published posthumously in 1925), and Thomas Mann's *Magic Mountain* (1924).

Various influences that characterized much of the writing from the 1920s were at work in these writers. An interest in the unconscious and the irrational was reflected in their work and that of others of about this time. Two important sources of this influence were Friedrich Nietzsche, a German philosopher to whom both Gide and Mann, for example, were much indebted, and Sigmund Freud, whose psychoanalytical works, by the 1920s, had had a telling influence on Western intellectuals. A shift away from 19th-century assumptions and styles was not limited to writers of fiction. André Breton's first *Manifeste du surréalisme* (1924; "Manifesto of Surrealism") was the first formal statement of a movement that called for spontaneity and a complete rupture with tradition. Surrealism showed the influence of Freud in its emphasis on dreams, automatic writing, and other antilogical methods and, although short-lived as a formal movement, had a lasting effect on much 20th-century art and poetry. The uncertainty of the new age and the variety of attempts to deal with it and give it

The influence of Nietzsche and Freud

some artistic coherence can be seen also in Rainer Maria Rilke's *Duino Elegies* and *Sonnets to Orpheus* (1923); in T.S. Eliot's *Waste Land* (1922); and in Luigi Pirandello's play about the instability of identity, *Henry IV* (1922).

The international and experimental period of Western literature in the 1910s and 1920s was important not only for the great works it produced but also because it set a pattern for the future. What was clearly revealed in the major works of the period was an increasing sense of crisis and urgency, doubts as to the 19th century's faith in the psychological stability of the individual personality, and a deep questioning of all philosophical or religious solutions to human problems. In the 1930s these qualities of 20th-century thought were not abandoned but, rather, were expanded into a political context, as writers divided into those supporting political commitment in their writing and those reacting conservatively against such a domination of art by politics. Nor did World War II resolve the debate concerning political commitment—issues similar to those that exercised major creative imaginations of the 1930s were still very much alive during the last quarter of the century.

The scarcity of great writers after World War II

It would be tempting to explain what seemed to be a relative scarcity of great writers in the period after World War II as an inevitable result of the cumulative pressure of disturbing social and technological developments accelerated by that war. Under such fluctuating and doubtful circumstances, it would not seem altogether strange if writing and reading, as traditionally understood, should cease. Indeed, in certain technologically highly developed countries, such as the United States, the printed word itself seemed to some critics to have lost its central position, having been displaced in the popular mind by a visual and aural electronic culture that did not need the active intellectual participation of its audience. Thus the communications media that helped to create something resembling an international popular culture in many Western countries did nothing to make the question of literary value easier to answer. Given the extraordinary conditions in which a modern writer works, it was not surprising that reputations were difficult to judge, that radical experimentation characterized many fields of literature, and that traditional forms of writing were losing their definition and were tending to dissolve into one another. Novels might acquire many features of poetry or be transformed into a kind of heightened nonfictional reportage, while experimentation with typography gave poems an appearance of verbal paintings, and dramatic works, shorn of anything resembling a traditional plot, became a series of carefully orchestrated gestures or events. But formal experimentation was only part of the picture, and to say that modern writing since World War II has been primarily experimental would be to ignore other characteristics that writing acquired earlier in the century and that still continued to be issues. Most good critics felt that there was no lack of good literature being written, despite the lack of major reputations and despite the possibly transitional nature of much of the period's work in its variety of styles and subjects.

BIBLIOGRAPHY

General introductory works: FERNAND BALDENSBERGER and WERNER P. FRIEDERICH, *Bibliography of Comparative Literature* (1950, reprinted 1960), comprehensive coverage of literatures, classified by author, country, and according to genre, theme, etc.; ANTONY BRETT-JAMES, *The Triple Stream: Four Centuries of English, French, and German Literature, 1531–1930* (1953, reprinted 1977), contains useful comparative tables arranged chronologically; JOHN M. COHEN, *A History of Western Literature*, rev. ed. (1963), a useful factual and critical account; JEAN-ALBERT BEDE and WILLIAM B. EDGERTON (gen. eds.), *Columbia Dictionary of Modern European Literature*, 2nd ed. (1980), contains over 1,800 articles, including surveys, critical essays, and biographies, with bibliographies; WILLIAM F. THRALL and ADDISON HIBBARD, *A Handbook to Literature*, 4th ed. by C. HUGH HOLMAN (1980), good coverage of literary terms and movements; RENE WELLEK, *A History of Modern Criticism: 1750–1950*, 4 vol. (1955–), surveys theory in the major European nations since the Renaissance.

Ancient period: JULIUS A. BEWER, *The Literature of the Old Testament in Its Historical Development*, 3rd ed. (1962); SIR

PAUL HARVEY (ed.), *The Oxford Companion to Classical Literature* (1937, reprinted with corrections 1969), a basic reference work; MAURICE PLATNAUER (ed.), *Fifty Years of Classical Scholarship* (1954), contains valuable bibliographies; C.M. BOWRA, *Landmarks in Greek Literature* (1966); ROBERT FLACELIERE, *Histoire littéraire de la Grèce* (1962; Eng. trans., *A Literary History of Greece*, 1964), a systematic, readable survey of ancient literature; KENNETH J. DOVER et al., *Ancient Greek Literature* (1980); HERBERT J. ROSE, *A Handbook of Greek Literature from Homer to the Age of Lucian*, 4th ed. rev. (1961); and *A Handbook of Latin Literature from the Earliest Times to the Death of St. Augustine*, 3rd ed. (1954, reprinted with suppl. bibliog. 1966), two standard works; MICHAEL GRANT, *Roman Literature* (1954), a good short account; WILLIAM A. LAIDLAW, *Latin Literature* (1951), a useful introduction; JOHN WIGHT DUFF, *A Literary History of Rome*, 2nd–3rd ed., 2 vol., ed. by A.M. DUFF (1960–64, 3rd ed. reprinted 1979), standard introduction containing comprehensive and scholarly surveys, with supplementary bibliographies.

The Middle Ages: H.M. CHADWICK, *The Heroic Age* (1912, reprinted 1974), a comparative study covering Greek, Teutonic, Slavonic, and Celtic literatures and cultures; E.K. CHAMBERS, *The Mediaeval Stage*, 2 vol. (1903), a standard reference work; W.T.H. JACKSON, *The Literature of the Middle Ages* (1960), and (ed.), *The Interpretation of Medieval Lyric Poetry* (1980); CHARLES W. JONES (ed.), *Medieval Literature in Translation* (1950), a good, readable anthology; C.S. LEWIS, *The Allegory of Love* (1936, reprinted 1958), a critical study of courtly love in medieval literature; ROGER S. LOOMIS (ed.), *Arthurian Literature in the Middle Ages* (1959), an excellent survey.

The Renaissance: ERNST CASSIRER, PAUL O. KRISTELLER, and JOHN H. RANDALL (eds.), *The Renaissance Philosophy of Man* (1948), a collection of excerpts from the writings of Renaissance philosophers; JOSEPH A. MAZZEO, *Renaissance and Revolution: The Remaking of European Thought* (1966), useful background reading; HENRY O. TAYLOR, *Thought and Expression in the Sixteenth Century*, 2nd rev. ed., 2 vol. (1959), a good standard work.

The 17th century: MARIO PRAZ, *Studi sul concettismo*, 2 vol. (1934–46; Eng. trans., *Studies in Seventeenth-Century Imagery*, 2 vol., 1939–48), a standard and readable work; JOEL E. SPINGARN (ed.), *Critical Essays of the Seventeenth Century*, 3 vol. (1909, reprinted 1957), a valuable collection.

The 18th century: LILIAN R. FURST, *Romanticism in Perspective*, 2nd ed. (1979), on the Romantic movement in France, England, and Germany, and *The Contours of European Romanticism* (1979); JOHN B. HALSTED (ed.), *Romanticism* (1969), a collection of extracts and key essays by the leading figures of western European Romanticism, with a long introduction and chronological table; PAUL HAZARD, *La Pensée européenne au XVIII^e siècle* (1963; Eng. trans., *European Thought in the Eighteenth Century*, 1963); DANIEL MORNET, *The Development of Literature and Culture in the XVIIIth Century* (1954), an extensive survey; MARIO PRAZ, *La carne, la morte e il diavolo nella letteratura romantica* (1930; 5th ed., 1976; Eng. trans., *The Romantic Agony*, 2nd ed., 1951; reissued with corrections, 1970), a remarkable work, now a classic in this field; JOHN G. ROBERTSON, *Studies in the Genesis of Romantic Theory in the Eighteenth Century* (1923, reprinted 1962).

The 19th century: GEORG BRANDES, *Creative Spirits of the Nineteenth Century* (1923, reprinted 1967); BENEDETTO CROCE, *Poesia e non poesia*, 7th ed. (1964; Eng. trans., *European Literature in the Nineteenth Century*, 1924, reprinted 1967); JANKO LAVRIN, *Aspects of Modernism: From Wilde to Pirandello* (1935, reprinted 1968); JOHN LUCAS (ed.), *Literature and Politics in the Nineteenth Century* (1971).

The 20th century: ERICH AUERBACH, *Mimesis* (1946; Eng. trans. 1953), traces the development of Realism from ancient to modern times; CLEANTH BROOKS and ROBERT PENN WARREN (eds.), *Understanding Poetry*, 4th ed. (1976), and *Understanding Fiction*, 3rd ed. (1979); and CLEANTH BROOKS and ROBERT B. HEILMAN, *Understanding Drama*, rev. ed. (1948), detailed anthologies, with explications of texts—invaluable introductions, especially to the study of modern writing, though not limited to 20th-century literature; GYORGY LUKACS, *Studies in European Realism* (1950, reprinted 1972), an important work by a major Marxist critic; RENE WELLEK and AUSTIN WARREN, *Theory of Literature*, 3rd ed. rev. (1970), a general and theoretical work, especially useful as a guide to modern literature; CLAUDE MAURIAU, *L'Alittérature contemporaine* (1958; Eng. trans., *The New Literature*, 1959); STEPHEN SPENDER, *The Struggle of the Modern* (1963) and *The Thirties and After* (1978); EDMUND WILSON, *Avant-Garde* (1931), *The Triple Thinkers*, rev. ed. (1948, reprinted 1977), and *The Twenties: From Notebooks and Diaries of the Period*, ed. by LEON EDEL (1975).

(Ed.)

Locke

The English philosopher John Locke was an initiator of the Enlightenment in England and France, an inspirer of the U.S. Constitution, and the author of, among other works, *An Essay Concerning Human Understanding*, his account of human knowledge, including the “new science” of his day—*i.e.*, modern science.

By courtesy of the Governing Body of Christ Church, Oxford



Locke, oil painting by Sir Godfrey Kneller. In Christ Church, Oxford.

THE LIFE OF JOHN LOCKE

Early years. Locke was born in Wrington, Somerset, on Aug. 29, 1632, and reared in Pensford, six miles south of Bristol. His family was Anglican with Puritan leanings. His father, a country attorney of modest means, fought on the Parliamentary side in the Civil War—a fact that later helped him to find a place for his son in Westminster School, then controlled by a Parliamentary committee (though its headmaster, Richard Busby, was a Royalist). The training there was thorough, but Locke later complained of the severity of its discipline. In 1652 he entered Christ Church, Oxford. Puritan reforms at Oxford had not yet altered the traditional Scholastic curriculum of rhetoric, grammar, moral philosophy, geometry, and Greek; Locke found the course insipid and interested himself in studies outside the traditional program, particularly experimental science and medicine. He was graduated with a B.A. degree in 1656 and an M.A. two years later, around which time he was elected a student (the equivalent of fellow) of Christ Church. In 1660, as a newly appointed tutor in his college, Locke enthusiastically welcomed the end of the Puritan Commonwealth and the restoration of Charles II to the throne.

In 1661 Locke inherited a portion of his father's estate, which ensured a modest annual income. His studentship would eventually be subject to termination unless he took holy orders, which he declined to do. Not wishing to make teaching his permanent vocation, he taught undergraduates for four years only. He served as secretary to a diplomatic mission to Brandenburg in 1665, and on his return he was immediately offered, but refused, another diplomatic post. His papers of this period, his correspondence, and his commonplace books all testify to his chief interests at the time, *viz.*, natural science, on the one hand, and the study of the underlying principles of moral, social, and political life, on the other. To remedy the narrowness of his education he read contemporary philosophy, particularly that of René Descartes, the father of modern philosophy. But more than all, experimental science engaged his interest.

He collaborated with Robert Boyle, one of the founders of modern chemistry, who was a close friend, and, toward the end of the period, with another friend, Thomas Sydenham, an eminent medical scientist.

Association with Shaftesbury. It was as a physician that Locke first came to the notice of the statesman Lord Ashley (later to become the 1st earl of Shaftesbury). On a visit to Oxford in the summer of 1666, Lord Ashley required some medical attention and was introduced to Locke by a mutual acquaintance; the two immediately became friends. A royal mandate of that November secured Locke's studentship indefinitely. The following year, despite his having no medical degree and no desire to practice medicine, he joined Ashley's household at Exeter House in the Strand in London as family physician. He became Ashley's personal adviser not merely on medical matters but on his general affairs as well.

Ashley was a forceful, aggressive politician who had many enemies (some of them men of letters—for instance, Locke's schoolfellow, the poet laureate John Dryden). It is doubtful, however—if only in view of Locke's respect for him—whether Ashley was as evil as his enemies sometimes made him out to be. It is known that he stood firmly for a constitutional monarchy, for a Protestant succession, for civil liberty, for toleration in religion, for the rule of Parliament, and for the economic expansion of Britain; and that he continued to make this stand when many influential men were working against these aims. Since these were already aims to which Locke had dedicated himself, there existed from the first a perfect understanding between the statesman and his adviser, one that meant much to both. Ashley entrusted Locke with the task of negotiating his son's marriage with the daughter of the Earl of Rutland; he also made him secretary of the group that he had formed to increase trade with America, particularly with the southern colonies. Locke helped to draft a constitution for the new colony of Carolina, a document that extended freedom of worship to all colonists, denying admission only to atheists.

During the following decades, Locke persevered in his private studies, and many of his social meetings were in effect meetings with friends to discuss philosophical and scientific problems. As early as 1668 he had become a fellow of the newly formed (1663) Royal Society, which kept him in touch with scientific advances. It is known, too, that groups of friends (Lord Ashley; the physician John Mapletoft; Thomas Sydenham; Sydenham's physician colleague, James Tyrrell, who was also a divine; and others) met in his rooms, for one such meeting is mentioned in the preface of his *Essay Concerning Human Understanding*, in which he reports that, because of the difficulties that beset the participants, they resolved to devote their next meeting to discussing the powers of the mind in order, as they said, “to examine our own abilities and see what objects our understandings were, or were not, fitted to deal with.” Locke himself opened the discussion and, following the meeting, set out his view of human knowledge in two drafts (1671), still extant, which show the beginnings of the thinking that 19 years later would blossom into his famous *Essay*. In these London years, too, Locke encountered representatives of Cambridge Platonism, a school of Christian humanists, who, though sympathetic to empirical science, nonetheless opposed materialism because it failed to account for the rational element in human life. They tended to be liberal in both politics and religion. Insofar as they taught a Platonism that rested on belief in innately known Ideas, Locke could not follow them; but their tolerance, their emphasis on practical conduct as a part of the religious life, and their rejection of materialism were features that he found most attractive. This school was closely related in spirit to another school that

Student
years

Private
studies
and
discussions

influenced Locke at this time, viz., that of latitudinarianism. For the latter school, if a man confessed Christ, that alone should be enough to entitle him to membership in the Christian Church; conformity in nonessentials should not be demanded. These movements prepared Locke for the antidogmatic, liberal school of theology that he would later encounter in Holland, a school in revolt against the narrowness of traditional Calvinism.

In 1672 Ashley was raised to the peerage as the 1st earl of Shaftesbury and at the end of that year was appointed lord high chancellor of England. Though he soon lost favour and was dismissed, he did, while in office, establish the Council of Trade and Plantations, of which Locke was secretary for two years. Locke, however, who suffered greatly from asthma, found the London air and his heavy duties unhealthy, and in 1675 he had to return to Oxford.

Intellectual
contacts in
France

Six months later he departed for France, where he stayed for four years (1675–79), spending most of his time in Paris and Montpellier. In France during the 1670s, Locke made contacts that deeply influenced his view of metaphysics and epistemology, viz., with the Gassendist school and, particularly, with its leader, François Bernier. Pierre Gassendi, a philosopher and scientist, had rejected over-speculative elements in Descartes's philosophy and had advocated a return to Epicurean doctrines—i.e., to empiricism (stressing sense experience), to hedonism (holding pleasure to be the good), and to corpuscular physics (according to which reality consists of atomic particles). Knowledge of the external world, Gassendi held, depends upon the senses, though it is through reasoning that man may derive much further information from empirically gained evidence.

Upon Locke's return to England, he found the country torn by dissension. The heir to the throne, James (the brother of Charles II), was a Roman Catholic, whom the Protestant majority led by Shaftesbury wished to exclude from the succession. For a year Shaftesbury had been imprisoned in the Tower, but by the time Locke returned he was back in favour once more as lord president of the Privy Council. When he failed, however, to reconcile the interests of the King and Parliament, he was dismissed; in 1681 he was arrested, tried, and finally acquitted by a London jury. A year later he fled to Holland, where, in 1683, he died.

Later life. No one of Shaftesbury's known friends was now safe in Great Britain. Locke himself, who was being closely watched, crossed to Holland in September 1683.

Exile in Holland. Locke's sojourn in Holland was happier than he had expected it to be: his health improved, he made many new friends, and he found the leisure that enabled him to bring his thoughts on many subjects to fruition. Locke spent his first winter in Amsterdam and soon became friendly with a distinguished Arminian theologian, Philip van Limborch, pastor of the Remonstrants' church there—a friendship that lasted until Locke's death. The companionship of Philip and other friends made it easier to bear bad news from home: at Charles II's express command, Locke (in 1684) was deprived of his studentship at Christ Church. The next year his name appeared on a list sent to The Hague that named 84 traitors wanted by the English government. Locke went into hiding for a while but soon was able to move freely over Holland and became familiar with its different provinces.

Return to England and retirement to Oates. Locke remained abroad for more than five years, until James II, who had become king in 1685, was overthrown. In the autumn of 1688, after it was announced that James had been presented with a male heir (and thus a Roman Catholic successor), the King's opponents invited his Protestant nephew and son-in-law, William of Orange in the Netherlands, to seize the throne. The King offered little resistance. Locke himself in February 1689 crossed in the party that accompanied the Princess of Orange, now to be crowned Queen Mary II of England. The triumph was complete; Locke was home again, although not without a nostalgia for the Holland that he had come to love. He now took little part in public life. He refused ambassadorial posts but accepted a membership in the Commission of Appeals. (Much later, in 1696, he

was appointed a commissioner in the resuscitated Board of Trade and Plantations, however, and for four years played a leading part in its deliberations.) But the London air again bothered him, and he was forced to leave the city for long visits to his friends in the country. In 1691 he retired to Oates, the house of his friends Sir Francis and Lady Masham in Essex, and subsequently made only occasional visits to London. Nonetheless, he was not without influence in these last years of his life, for he was the intellectual leader of the Whigs. Their principal parliamentarians were frequently old friends of Locke, and the younger generation—particularly the ablest of them all, John Somers, who soon became lord chancellor—turned to him constantly for guidance. In "the glorious, bloodless revolution," the main aims for which Shaftesbury and Locke had fought were achieved—even though in William's reign strong Tory pressures limited the extent of the reform. First and foremost, England became a constitutional monarchy, controlled by Parliament. Second, real advances were made in securing the liberty of subjects in the law courts, in achieving a greater (though far from complete) measure of religious toleration, and in assuring freedom of thought and expression. Locke himself drafted the arguments that his friend Edward Clarke used in the House of Commons in arguing for the repeal of the restrictive Act for the Regulation of Printing. The act was abolished in 1695 and the freedom of the press was secured.

Influence
during
retirement

Publication of his works. The main task of this last period of his life, however, was the publication of his works, which had been the product of long years of gestation. The *Epistola de Tolerantia (A Letter Concerning Toleration)*, 1689) was published anonymously at Gouda in 1689. Locke had been reflecting on this topic from his early days at Oxford. Though his correspondence and a paper that he wrote in 1667 show his support for toleration in religion, in 1660–61 he wrote two tracts on this theme (not published until 1967) that are surprisingly conservative. *Two Treatises of Government* (1690) was also the fruit of years of reflection upon the true principles in politics, a reflection resting on Locke's own observations. In all of these social and political issues, Locke saw that the ultimate factor is man's nature. To understand man, however, it is not enough to observe his actions: one must also inquire about his capacities for knowledge. Locke had been conscious of this point in writing his paper on the "Law of Nature" as early as 1663. In 1671, as has been seen, he set out to write a book about human knowledge. *An Essay Concerning Human Understanding*, which was not published, however, until December 1689 (all copies dated 1690)—nor was it wholly completed even then, for Locke made changes, sometimes substantial ones, in three of the four following editions. (See below, *Locke's philosophy*.)

Last years. Locke's last years were spent in the peaceful retreat of Oates. His hostess was a woman with whom he had been acquainted for many years—Lady Masham, or Damaris, the daughter of Ralph Cudworth, one of the Cambridge Platonists, by whom Locke had been significantly influenced. He found friendship and comfort in this household. Many of his friends visited him there: Sir Isaac Newton, who came to discuss the Epistles of St. Paul, a subject of great interest to both; his nephew and heir Peter King, destined to become lord high chancellor of England; and Edward Clarke with his wife and children, for whom Locke had great affection. Locke had written a series of letters to Edward Clarke from Holland, advising him on the best upbringing for his son. These letters formed the basis of his influential *Some Thoughts Concerning Education* (1693), setting forth new ideals in that field. He wrote and published pamphlets on matters of economic interest, on rates of interest, on the coinage of the realm, and, more widely, on trade (defending mercantilist views). In 1695 he published a dignified plea for a less dogmatic Christianity in *The Reasonableness of Christianity*.

Later
writings

John Locke died on Oct. 28, 1704, and was buried in the parish church of High Laver. "His death," wrote Lady Masham, "was like his life, truly pious, yet natural, easy and unaffected." This account of his character by one

who knew him well seems singularly appropriate. He was orderly, careful about money, occasionally parsimonious, abstemious, and, though naturally emotional and hot-tempered, controlled and disciplined. He had a great love of children, and friendship was for him a necessity. Both in his books and in his life are found the marks of the prudence and wisdom for which he was famed.

LOCKE'S PHILOSOPHY

Theory of knowledge. Locke was thoroughly suspicious of the view that a thinker could work out by reason alone the truth about the universe. Much as he admired Descartes, he feared this speculative spirit in him, and he despised it in the Scholastic philosophers. In this sense he rejected metaphysics. Knowledge of the world could only be gained by experience and reflection on experience, and this knowledge was being gained by Boyle, Sydenham, Christiaan Huygens, and Newton. They were the true philosophers who were advancing knowledge. Locke set himself the humbler task, as he conceived it, of understanding how this knowledge was gained. What was "the original, certainty, and extent of human knowledge, together with the grounds and degrees of belief, opinion, and assent"?

Empiricism. As for "the original," the answer was plain. Knowledge of the world began in sense perception, and self-knowledge in introspection, or "reflection" in Locke's language. It did not begin in innate knowledge of maxims or general principles, and it did not proceed by syllogistic reasoning from such principles. In the 17th century there had been much vague talk about innate knowledge, and in Book I of his *Essay Concerning Human Understanding* Locke examines this talk and shows its worthlessness. In Book II of his *Essay* he begins by claiming that the sources of all knowledge are sense experience and reflection; these are not themselves, however, instances of knowledge in the strict sense, but they provide the mind with the material of knowledge. Locke calls the material so provided "ideas." Ideas are objects "before the mind," in the sense not that they are physical objects but that they represent them. Locke distinguishes ideas that represent actual qualities of objects (such as size, shape, or weight) from ideas that represent perceived qualities, which do not exist in objects except as they affect observers (such as colour, taste, or smell). Locke designates the former primary qualities and the latter secondary qualities.

Locke proceeds to group and classify the ideas, with a view to showing that the origin of all of them lies in sensation and reflection. Although ideas are immediately "before the mind," not all of them are simple. Many of them are compounded, and their simple parts can be revealed on analysis. It is these simple ideas alone that are given in sensation and reflection. Out of them the mind forms complex ideas, though Locke is ambiguous on this point. For while he uses the language of "forming" or "compounding" and speaks of the "workmanship" of the mind, the compounding is frequently in accordance with what is perceived "to go together" and is not arbitrary.

Locke's reflections upon cause and effect, had they been elaborated, would undoubtedly have led him into acute difficulties. He does admit one failure. As an empiricist he can give no account of the idea of substance; it is, he thinks, essential and not to be denied, and yet it is not a simple idea given in sensation or reflection nor is it derived from simple ideas so given. In fact he can say little of it; it is "a-something-I-know-not-what." Thus, the case for empiricism cannot be said to be entirely established by Book II, but Locke thinks it strong enough for him to persist in the view that knowledge of the physical world is wholly derived from sense perception.

Self-knowledge. Some ideas are not of things outside the mind but are reflexive and internal. Locke finds it necessary to classify these in Book II and in doing so sets down the foundations of empirical psychology. His source of information is introspection and rarely the observation of behaviour. His account of sense perception is celebrated for its appreciation of the part that the interpretative mind plays in perceiving, and some of his farsighted observations on the relations between the senses, particularly vision

and touch, have profoundly affected subsequent thought. He makes valuable remarks on memory, on discerning, on comparing, on madness, on pleasure and pain, on the emotions, and on the association of ideas.

Locke holds that man has an intuitive knowledge of his own existence and supposes that man exists as material and immaterial substance, but he is none too clear about this and at one point plays with the idea that man is simply material substance to which God has "superadded" a power of thinking. Locke's most valuable contribution, however, is his account of personal identity. Having distinguished between different types of identity, he argues that personal identity depends on self-consciousness (that is, I am the person who did so-and-so 20 years ago because I can remember myself doing it).

Language. According to Locke, Book III on language "cost [him] more pains" than any other book of his *Essay*; yet it is the book that has been most neglected. To understand thinking and knowing one must understand language as the means of thought and communication. Words are conventional signs; however, according to Locke, signs do not directly represent things but rather ideas of things. Thus, Locke carries a theory of ideas into his account of language. Frequently, the idea signified by the word is not clear, and sometimes words are used even when there are no ideas corresponding to them. This is particularly so in the case of general words, without which language would be so impoverished as to lose most of its worth. The use of general words, in Locke's mind, is bound up with the theory of universals. Does the general word stand for a particular idea that is used in a representative capacity? Or is the universal nothing more than a creation of the mind, through abstraction, to which is attached a name? In considering natural substances, Locke is inclined strongly toward a conceptualism according to which the use of general words is possible only because they signify "nominal essences." In this view what is meant is not the real essence but an abstract concept, something brought about through the "workmanship of the understanding." Locke also discusses the names of simple ideas and of relations, and it is interesting to find the crude beginnings of a discussion of what were later to be called logical or operative words. Book III contains also a valuable account of definition, which denies the theory that all definition must be *per genus et differentiam* (by comparison and contrast). The final chapters deal with the inevitable imperfections of language and with avoidable abuses.

Conclusions. In Book IV, Locke discusses the nature and extent of human knowledge. The tone is more rationalistic than that of the previous books because the skepticism that emanated from his empiricism drove him to find the ideal of knowledge in the indubitable certainties of mathematics. There he was on common ground with the rationalists of his day, and indeed the direct influence of Descartes seems to be observable in the opening chapters of Book IV. Knowledge is perception, not sense perception but intellectual perception or intuition, frequently gained by a deliberate process of demonstration. But, even when this is so, each step in the demonstration is observed intuitively, so that knowledge in the strict sense is essentially intuitive.

Unfortunately, what can be intuited and demonstrated is limited. Strict knowledge is not confined entirely to mathematics, but the intuition of relations within the physical world is impossible. Books II and III have shown that ideas and nominal essences can be grasped directly and that the inner nature of real things cannot be known, so that "science," in the exact sense of perfectly certain knowledge, is not possible in this sphere. The only possibility of intuiting is that within the world of ideas, an ideal world that is for Locke empirically derived and not intellectual in character. Knowledge in general terms he accordingly defines as the intuition or "perception of the connection of and agreement, or disagreement and repugnancy, of any of our ideas." Within the realm of ideas indubitable knowledge can be gained, but when dealing with ideas "whose archetypes are without them" the position is uncertain.

In spite of this, and somewhat inconsistently, Locke

General words

Simple and complex ideas

Knowledge and intuition

thinks that knowledge approaches certainty in the "sensitive knowledge" of the existence of physical things. Further, knowledge of one's own existence is intuited. In these cases knowledge that is not an apprehension of a relation between ideas is nonetheless certain. But Locke makes it clear that, for the most part, knowledge of the physical world or of oneself is probable and rests not on intuition but on judgment; it is assenting to a proposition on the strength of the evidence, and there may be degrees of assent and wrong assent or error. Locke recognizes the need for a logic of probability, though he does little himself to meet that need. Yet it should be added that the important regular-sequence theory of induction, afterward developed by George Berkeley and David Hume, is put forward in the pages of Locke's *Essay*.

Political theory. Locke's most important work on political philosophy is that entitled *Two Treatises of Government*. The first treatise is a refutation of Sir Robert Filmer's *Patriarcha*, a defense of the divine right of kings that was written in the mid-17th century; the second and more important treatise refutes the absolutist theory of government as such.

Locke defines political power as "a right of making laws, with penalties of death, and consequently all less penalties for the regulating and preserving of property and of employing the force of the community in the execution of such laws, and in the defence of the commonwealth from foreign injury, and all this only for the public good." Government is thus a trust, forfeited by a ruler who fails to secure the public good. The ruler's authority, that is to say, is conditional rather than absolute. Nor does the individual surrender all his rights when he enters a civil society. He has established his right to property by "mixing his labour" with things originally given to mankind in common but now made his own by his labour. (Here in germ is the labour theory of value.) He has the right to expect political power to be used to preserve his property, in his own person and in his possessions, and the right to freedom of thought, speech, and worship. In fact the one right that he gives up in entering a civil society is the right to judge and punish his fellow man, which is his right in the state of nature. He quits his "executive power of the law of Nature" and "resigns it to the public"; he himself makes himself subject to the civil law and finds his freedom in voluntary obedience. To secure this freedom, Locke favoured a mixed constitution—the legislative should be an elected body, whereas the executive is usually a single person, the monarch—and he argues for a separation of legislative and executive powers. The people are ultimately sovereign, although it is not always clear in Locke's theory where the immediate sovereignty lies. But the people always have the right to withdraw their support and overthrow the government if it fails to fulfill their trust.

Moral philosophy. One searches in vain for a consistent moral theory in Locke. His view that morality can be a science, as certain as mathematics, is well known. This might imply a rationalism, and there are indeed rationalist trends in his moral philosophy—although sometimes when advocating a science of morals he seems to have in mind simply the possibility of an exact analysis of the terms used in moral discourse and the clarification of moral statements. At other times, he puts forward a hedonist theory. "That we call *good* which is apt to cause or increase pleasure or diminish pain in us." But not every good is moral good: "Moral good and evil is only the conformity or disagreement of our voluntary actions to some law, whereby good or evil is drawn on us, from the will and power of the law-maker." In this view law rests on God's will, "the true ground of morality," though in saying this Locke does not appear to be consistent with what he says elsewhere of the law of nature.

Theory of education. A good education, as set forth by Locke in *Some Thoughts Concerning Education*, attends to both the physical and the mental. The body is not to be coddled; on the contrary, it is necessary that it should be hardened in various ways. The good educator insists on exercise, play, and plentiful sleep, "the great cordial of nature." Young children should be allowed to give

vent to their feelings and should be restrained rarely. As for mental training, character comes first before learning; the educator's aim is to instill virtue, wisdom, and good breeding into the mind of the young. Parents, too, must interest themselves in their children's upbringing and, as far as possible, have them near; for no educative force is more powerful than the good example of parents. A stock of useful knowledge must be imparted: modern languages and Latin; geography and history; mathematics, as "the powers of abstraction develop"; and later civil law, philosophy, and natural science. For recreation, training in the arts, crafts, and useful hobbies should be available.

Religion. Locke's reaction against the "enthusiasm" of the sects in his youth had been sharp, and he disliked religious fanaticism throughout his life. He was a broad, tolerant Anglican anxious to heal the breach in English Protestant ranks. His own views on church government and on the priesthood were close to those of the dissenters, and he favoured the liberal views of the latitudinarians, of the Cambridge Platonists, and of the Remonstrants of Holland. This becomes manifest in *The Reasonableness of Christianity*. Two essentials, and two alone, he thinks, are involved in being a Christian: first, that a man should accept Christ as God's Messiah and, second, that he should live in accordance with Christ's teaching. His point of view is not far removed from that of the Deists on the one hand and the Unitarians on the other, yet he cannot be grouped with them. Christianity, though reasonable, needs revelation as well as reason, for human reason alone is inadequate: there is an experience of God "through His Spirit" without which all religion is empty. However, any act of persecution in the name of religious truth is wholly unjustified, since our knowledge and understanding are so confined. Each individual is a moral being, responsible before God, and this presupposes freedom. By the same token, no compulsion that is contrary to the will of the individual can secure more than an outward conformity.

Influence. Locke's faith in the salutary, ennobling powers of knowledge justifies his reputation as the first philosopher of the Enlightenment. In a broader context, he founded a tradition of thought that would span three centuries, in the schools of British empiricism and American pragmatism. In developing the Whig ideology underlying the Exclusion Controversy and the Revolution of 1688, Locke formulated the classic expression of liberalism, which was to inspire both the shapers of the American Revolution and the authors of the U.S. Constitution. Locke's influence remained strongly felt in the West in the 20th century, as notions of mind, freedom, and authority continued to be challenged and explored.

MAJOR WORKS

PHILOSOPHY, RELIGION, AND EDUCATION: *An Essay Concerning Humane [sic] Understanding* (1690); *Epistola de Tolerantia* (1689); *A Letter Concerning Toleration*, trans. by William Popple, 1689; *A Second Letter Concerning Toleration* (1690); *A Third Letter for Toleration* (1692); *Some Thoughts Concerning Education* (1693); *The Reasonableness of Christianity, as Delivered in the Scriptures* (1695); *A Vindication of the Reasonableness of Christianity* (1695); *A Second Vindication of the Reasonableness of Christianity* (1697); *Of the Conduct of the Understanding*, in *Posthumous Works of Mr. John Locke* (1706).

POLITICAL PHILOSOPHY AND ECONOMICS: *Two Treatises of Government* (1690); *Some Considerations of the Consequences of the Lowering of Interest, and Raising the Value of Money* (1692); *Short Observations on a Printed Paper, Intituled, for Encouraging the Coining Silver Money in England, and After for Keeping It Here* (1695); *Further Considerations Concerning Raising the Value of Money* (1695).

RECOMMENDED EDITIONS: A complete edition is *The Works of John Locke*, new ed. corrected, 10 vol. (1823, reprinted 1963). There is no complete modern edition of Locke's works, although several volumes have appeared in the Oxford Press series, "The Clarendon Edition of the Works of John Locke"; the first of these was a critical edition of *An Essay Concerning Human Understanding*, edited by Peter H. Niddich (1975, reprinted 1979). Useful editions of other single works include *A Letter Concerning Toleration*, edited by James Tully (1983); *Two Treatises of Government*, edited by Peter Laslett, 2nd ed. (1967, reprinted 1970), a critical edition; and *The Educational Writings of John Locke: A Critical Edition with Introduction and Notes*, edited by James L. Axtell (1968).

Freedom through law

Reason and revelation

BIBLIOGRAPHY

Biographies: Early works include LORD KING, *The Life and Letters of John Locke*, new ed. (1858, reissued 1984), an amateurish work but based on the Lovelace Collection of Locke papers in the possession of Peter King's family; and H.R. FOX BOURNE, *The Life of John Locke*, 2 vol. (1876, reprinted 1969), a detailed study, based on secondary sources. MAURICE W. CRANSTON, *John Locke: A Biography* (1957, reissued 1985), is now the standard biography. An outstanding resource is E.S. DE BEER (ed.), *The Correspondence of John Locke* (1976–), part of "The Clarendon Edition of the Works of John Locke"; 7 of 8 vol. have appeared to 1986.

Commentaries: JOHN W. YOLTON, *Locke: An Introduction* (1985); and JOHN DUNN, *Locke* (1984), provide general accounts of Locke's life and work. For Locke's theory of knowledge, see R.S. WOOLHOUSE, *Locke* (1983); and JAMES GIBSON, *Locke's Theory of Knowledge and Its Historical Relations* (1917, reprinted 1968), another useful introductory essay, if somewhat old-fashioned in its approach. For a survey of Locke's thought, see RICHARD I. AARON, *John Locke*, 3rd ed. (1971, reprinted 1973); D.J. O'CONNOR, *John Locke* (1952, reissued 1967); and JOHN W. YOLTON, *John Locke and the Way of Ideas* (1956, reprinted 1968), a study based on Locke's unpublished as well as his published writings.

Specialized commentaries on Locke's epistemology are found in JOHN W. YOLTON, *Locke and the Compass of Human Understanding: A Selective Commentary on the "Essay"* (1970); J.L. MACKIE, *Problems from Locke* (1976); and I.C. TIPTON (ed.), *Locke on Human Understanding: Selected Essays* (1977). Political theory is covered in STERLING POWER LAMPRECHT, *The Moral and Political Philosophy of John Locke* (1918, reprinted 1962); GERAINT PARRY, *John Locke* (1978); J.W. GOUGH, *John Locke's Political Philosophy: Eight Studies*, 2nd ed. (1973); and M. SELIGER, *The Liberal Politics of John Locke* (1968), an exposition and a defense of Locke's arguments for political freedom. W. VON LEYDEN, *Hobbes and Locke: The Politics of Freedom and Obligation* (1981); RICHARD H. COX, *Locke on*

War and Peace (1960, reprinted 1982); and C.B. MACPHERSON, *The Political Theory of Possessive Individualism: Hobbes to Locke* (1962, reprinted 1983), explore the relationship between Locke's political thought and that of Thomas Hobbes. See also JOHN DUNN, *The Political Thought of John Locke: An Historical Account of the Argument of the "Two Treatises of Government"* (1969, reprinted 1982), a survey of Locke's thought in the context of his intellectual environment; and RAYMOND POLIN, *La Politique morale de John Locke* (1960, reprinted 1984), on Locke's liberalism from the perspective of a French historian of ideas. JAMES TULLY, *A Discourse on Property: John Locke and His Adversaries* (1980, reissued 1982); GORDON J. SCHOCHET, *Life, Liberty and Property: Essays on Locke's Political Ideas* (1971); and J.G.A. POCOCK and RICHARD ASHCRAFT, *John Locke* (1980), discuss Locke's defense of the natural right to property. See also KAREN IVERSEN VAUGHN, *John Locke: Economist and Social Scientist* (1980), for Locke's ideas on economics; and KENNETH DEWHURST, *John Locke, 1631–1704, Physician and Philosopher* (1963, reprinted 1984), on his career as a practitioner and theorist of medical science. Research in progress, queries, and corrections to published work on Locke are reported in *The Locke Newsletter* (annual).

Bibliographies: H.O. CHRISTOPHERSEN, *A Bibliographical Introduction to the Study of John Locke* (1930, reprinted 1968), is still useful, although its references have been assimilated into a larger, more recent work, JEAN S. YOLTON and JOHN W. YOLTON, *John Locke: A Reference Guide* (1985)—both cover mainly secondary sources. JOHN C. ATTIG (comp.), *The Works of John Locke: A Comprehensive Bibliography from the Seventeenth Century to the Present* (1985), tracks the various editions and translations of Locke's writings and places them in historical context. See also ROLAND HALL and R.S. WOOLHOUSE, *80 Years of Locke Scholarship: A Bibliographical Guide* (1983); and P. LONG, *A Summary Catalogue of the Lovelace Collection of the Papers of John Locke in the Bodleian Library* (1959), a guide to the most important source of manuscript material.

(R.Aa./Ed.)

The History and Kinds of Logic

The major task of logic has been to establish a systematic way of deducing the logical consequences of a set of sentences. In order to accomplish this, it is necessary first to identify or characterize the logical consequences of a set of sentences. Then the procedures for deriving conclusions from a set of sentences need to be examined to verify that all logical consequences, and only those, are deducible from that set. Finally, in recent times, the question has been raised whether all the truths regarding some domain of interest can be contained in some specifiable deductive system.

In this article the several formal systems of modern logic

are examined, as is the general study of these systems themselves. Consideration is given to the application of logical investigations to its original area of concern, argumentative discourse, as well as to new areas such as epistemology and deontology (the logic of duties). The historical development of logic from ancient Greece through the medieval period to modern times is delineated.

(M.L.Sc.)

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, sections 10/11 and 10/12, and the *Index*.

This article is divided into the following sections:

- | | |
|--|--|
| Logic systems 230 | Precursors of ancient logic |
| Formal logic 230 | Aristotle |
| General observations | Theophrastus of Eresus |
| The propositional calculus | The Megarians and Stoics |
| The predicate calculus | Late representatives of ancient Greek logic |
| Modal logic | Medieval logic 265 |
| Set theory | Transmission of Greek logic to the Latin West |
| Metalogic 244 | Arabic logic |
| Nature, origins, and influences of metalogic | The revival of logic in Europe |
| Nature of a formal system and of its formal language | Developments in the 13th and early 14th centuries |
| Discoveries about formal mathematical systems | Late medieval logic |
| Discoveries about logical calculi | Modern logic 268 |
| Model theory | The 16th century |
| Applied logic 251 | The 17th century |
| The critique of forms of reasoning | Leibniz |
| Epistemic logic | The 18th and 19th centuries |
| Practical logic | 20th-century logic 275 |
| Logics of physical application | Russell and Whitehead's <i>Principia Mathematica</i> |
| Computer design and programming | 20th-century set theory |
| Hypothetical reasoning and counterfactual conditionals | Logic and philosophies of mathematics |
| The history of logic 261 | Logic narrowly construed |
| Origins of logic in the West 261 | Nonmathematical formal logic |
| | Bibliography 281 |

Arguments **General considerations.** From its very beginning, the field of logic has been occupied with arguments, in which certain statements, the premises, are asserted in order to support some other statement, the conclusion. If the premises are intended to provide conclusive support for the conclusion, the argument is a deductive one. If the premises are intended to support the conclusion only to a lesser degree, the argument is called inductive. A logically correct deductive argument is termed valid, while an acceptable inductive argument is called cogent. The notion of support is further elucidated by the observation that the truth of the premises of a valid deductive argument necessitates the truth of the conclusion: it is impossible for the premises to be true and the conclusion false. The truth of the premises of a cogent inductive argument, on the other hand, confers only a probability of truth on its conclusion: it is possible for the premises to be true while the conclusion is false.

Logic is not concerned to discover premises that persuade an audience to accept, or to believe, the conclusion. This is the subject of rhetoric. The notion of rational persuasion is sometimes used by logicians in the sense that, if one were to accept the premises of a valid deductive argument, it would not be rational to reject the conclusion; one would in effect be contradicting oneself in practice. The case of inductive logic will be considered below.

From the above characterization of arguments, it is evident that they are always advanced in some language, either a natural language such as English or Chinese or, possibly, a specialized technical language such as mathematics. To develop rules for determining the validity of deductive arguments, the statements comprising the argument must be analyzed in order to see how they relate to one another. The analysis of the logical forms of arguments can be accomplished most perspicuously if the statements of the argument are framed in some canonical form. Additionally, when stated in a regimented format, various ambiguities or other defects of the original statements can be avoided.

When they are stated in a natural language, some arguments appear to give support to their conclusions or to confute a thesis. Such a defective, although apparently correct, argument is called a fallacy. Some of these errors in argument occur often enough that types of such fallacies are given special names. For example, if one were to attack the premises of an argument by casting aspersions on the character of the proponent of the argument, this would be characterized as committing an *ad hominem* fallacy. The character of the proponent of an argument has no relevance to the validity of the argument. There are several other fallacies of relevance, such as threatening the audience (*argumentum ad baculum*) or appealing to their feelings of pity (*argumentum ad misericordiam*).

The other major grouping of fallacies concerns those apparently correct arguments whose plausibility depends on some ambiguity. For an argument to be valid it is required that the terms occurring in the argument retain one meaning throughout. Subtle shifts of meaning that destroy the correctness of any argument can occur in natural language expressions:

Today chain-smokers are rapidly disappearing.
Karen is a chain-smoker.
Therefore, today Karen is rapidly disappearing.

Clearly what is intended in the first premise is that the class of chain-smokers is becoming a smaller class, not that the individuals in the class are undergoing any change. A well-known, classic example of incorrect reasoning based on an ambiguity arising from the grammatical construction employed, the so-called amphiboly, is the case of Croesus, king of Lydia in the 6th century BC, who was considering invading Persia. When he consulted the oracle at Delphi, he is reported to have received the following reply: "If Croesus goes to war with Cyrus (the king of Persia), he will destroy a mighty kingdom." Croesus inferred that his campaign would be successful, but in fact he lost, and consequently his own mighty kingdom was destroyed.

One of the first and best-known—and most successful—attempts to provide a regimented framework within which

some important deductive arguments could be recognized as valid or invalid was that of Aristotle. Many arguments are composed of premises and conclusions that are stated or could be restated as categorical propositions. Categorical propositions may be distinguished first by their quality, either affirmative or negative. An affirmative categorical proposition asserts that all or some of a class of objects are included in another class of objects (e.g., "All whales are mammals"), while a negative categorical proposition asserts that all or some of a class of objects are not included in another class of objects (e.g., "Some pets are not dogs").

Secondly, categorical propositions may be distinguished by their quantity, either universal or particular. When the assertion is that all of a class of objects are or are not included in another class of objects, the proposition is universal. When only some (precisely, at least one) of a class are or are not included in another, the proposition is particular.

The two distinguishing features above lead to four types of categorical proposition:

A: universal affirmative	All <i>A</i> 's are <i>B</i> 's.
E: universal negative	No <i>A</i> 's are <i>B</i> 's.
I: particular affirmative	Some <i>A</i> 's are <i>B</i> 's.
O: particular negative	Some <i>A</i> 's are not <i>B</i> 's.

The letters to the left, A, E, I, and O, are the standard labels for these types of propositions. The expressions in the right column are schematic sentences, requiring, in this case, English phrases referring to classes of objects where *A* and *B* are located. Some examples of categorical propositions in this standard form are:

A:	All games are enjoyable activities.
E:	No wars are enjoyable activities.
I:	Some women are soldiers.
O:	Some women are not soldiers.

Not all arguments in ordinary contexts are expressed in categorical propositions. Indeed, most are not. The sample *A* proposition above would more likely be expressed as: "All games are enjoyable." But *enjoyable* is an adjective and does not refer to a class of objects. The adjective must be replaced by a noun phrase to obtain a proper categorical proposition. In all cases, propositions must be expressed using two noun phrases joined by the appropriate copula, a form of the verb *to be*.

Original: Some sailors are dancing.

Rewritten: Some sailors are persons who are dancing.

(Note that "Some sailors are dancers" is not quite right, since a dancer may not actually be dancing at the moment.)

Most languages contain many more verbs than the standard copula; hence, there are many grammatical statements that do not use variations of this verb. These sentences must be rewritten as well:

Original: All dogs bark.

Rewritten: All dogs are animals that bark.

Even variations of the verb *to be* must be rewritten:

Original: Some lucky person will win the lottery.

Rewritten: Some lucky persons are persons who will win the lottery.

Another difficulty with the requirement that all arguments be expressed using categorical propositions is that some arguments involve reference to one individual. The sentence "Socrates is a Greek" is considered to be a singular proposition. Some logicians allow such sentences in arguments and treat them as universal categorical propositions. It is usually better, however, to rewrite such sentences as explicit categorical propositions:

All persons identical to Socrates are Greeks.

The class referred to by the subject term "persons identical to Socrates" has one and only one object in it—namely, Socrates himself.

A natural language usually has various rhetorical devices for expressing quantifiers, and some languages—English, for example—occasionally do not even express the quantifier, letting the grammatical construction convey that

information instead. We find "A cow is a mammal" referring to cows in general, so it would be regimented as "All cows are mammals." Examples of noncategorical quantifiers along with appropriate translations into categorical propositions are:

- Original: A few scientists are dullards.
- Rewritten: Some scientists are dullards.
- Original: Not everyone who runs for office is elected.
- Rewritten: Some persons who run for office are not elected persons.
- Original: All entrants can't be winners.
- Rewritten: Some entrants are not winners.
- Original: Automobiles are not toys.
- Rewritten: No automobiles are toys.

Condi-
tional
sentences

Conditional sentences have the form "If . . . , then ____." If the antecedent ("if" clause) and the consequent ("then" clause) refer to the same class of objects, the conditional can be rewritten in categorical form. Otherwise, it cannot be rewritten and must be dealt with differently (see below *Other argument forms*). Some conditionals whose antecedent and consequent refer to the same class of objects are:

1. If an animal is a tiger, (then) it's a carnivore.
2. If it's a snake, then it's not a mammal.
3. A student will succeed if he or she studies assiduously. (Note the reversal of the clauses.)

These are rewritten in categorical form as:

1. All tigers are carnivores.
2. No snakes are mammals.
3. All students who study assiduously are students who will succeed.

When the antecedent and consequent refer to different classes, such rewriting is not possible (e.g., "If the president is reelected, then I shall never vote again").

Finally there are such locutions as "Only" (or "None but"), "The only," and "All except" (or "All but"). When it is asserted that only *A*'s are *B*'s, it is not claimed that *A*'s are *B*'s. Rather, it is claimed that, if anything is a *B*, then it is also an *A*. So, for example, if it is asserted that only entrants are prizewinners, no one is asserting that all entrants will win a prize. What is asserted is that all prizewinners are entrants. The case "The only" is quite different. Here, "The only winners are Texans" is expressed by the proposition "All winners are Texans." The phrase "All except" introduces an exceptive proposition. It requires two categorical propositions to state everything asserted by an exceptive proposition. The statement "All except crew members abandoned ship" asserts that everyone who was not a crew member abandoned ship and that no crew member abandoned ship. Thus, two categorical propositions are needed to express this exceptive proposition:

All non-crew members are persons who abandoned ship.

No crew members are persons who abandoned ship.

Immediate inference. The simplest possible arguments that can be constructed from categorical propositions are those with one premise and, of course, one conclusion. These are called immediate inferences. In order to characterize the valid arguments with one premise, it is necessary to consider various transformations of a categorical proposition. One transformation switches the subject and predicate terms of a proposition, resulting in a proposition called the converse of the original.

Transformations

Original	Converse
A: All <i>A</i> 's are <i>B</i> 's.	All <i>B</i> 's are <i>A</i> 's.
E: No <i>A</i> 's are <i>B</i> 's.	No <i>B</i> 's are <i>A</i> 's.
I: Some <i>A</i> 's are <i>B</i> 's.	Some <i>B</i> 's are <i>A</i> 's.
O: Some <i>A</i> 's are not <i>B</i> 's.	Some <i>B</i> 's are not <i>A</i> 's.

Only in the cases of E and I propositions can one immediately infer the converse. That is, only these inferences by conversion are correct:

- No snakes are birds.
- ∴ No birds are snakes.
- Some cats are pets.
- ∴ Some pets are cats.

The obverse of a proposition is a more complicated transformation. The quality of the proposition is changed from affirmative to negative (or from negative to affirmative), and the predicate term is replaced by its negation (frequently formed by prefixing "non-"). Thus, "All *A*'s are *B*'s" becomes "No *A*'s are non-*B*'s," and similarly for the other three categorical propositions. The obverse of any categorical proposition is logically equivalent to the original and hence may be immediately inferred from it:

- No snakes are birds.
- ∴ All snakes are non-birds.
- Some cats are pets.
- ∴ Some cats are not non-pets.
- All whales are mammals.
- ∴ No whales are non-mammals.
- Some dogs are not friendly animals.
- ∴ Some dogs are non-friendly animals.

The contrapositive of a categorical proposition is formed by converting the proposition (switching subject and predicate terms) and then negating both the subject and predicate. Only in the cases of A and O propositions can the contrapositive be inferred as a valid conclusion:

- All whales are mammals.
- ∴ All non-mammals are non-whales.
- Some pets are not cats.
- ∴ Some non-cats are not non-pets.

In the cases of E and I propositions, the contrapositive does not follow as a valid conclusion.

These immediate inferences are frequently employed to transform propositions in an argument into a form that enables the more complex argument to be analyzed.

Categorical syllogisms. The next more complex form of argument is one with two categorical propositions as premises and one categorical proposition as conclusion. When arguments of this type have exactly three terms occurring throughout the argument and when the predicate term of the conclusion occurs in the first premise and the subject term of the conclusion occurs in the second premise, the argument is called a categorical syllogism.

The pattern of the types of categorical propositions as they occur in a syllogism, frequently indicated by the appropriate letters (A, E, I, O), is called the mood of the syllogism. Thus, possible moods are AAA, AIO, EIO, and so on. Within a given mood, the terms can occur in various patterns. The pattern in which the terms S, M, and P (subject, middle, and predicate) are arranged is called the figure of the syllogism. For instance, in the first premise the predicate term of the conclusion may appear first as the subject of the premise or it may occur last as the predicate of the premise. This is also true for the subject term of the conclusion when it occurs in the second premise. There are four possibilities:

Mood and
figure

Figure 1	Figure 2	Figure 3	Figure 4
M—P	P—M	M—P	P—M
S—M	S—M	M—S	M—S
∴ S—P	∴ S—P	∴ S—P	∴ S—P

Thus a syllogism in the fourth figure, with mood AAA, is called AAA-4:

- All *P*'s are *M*'s.
- All *M*'s are *S*'s.
- ∴ All *S*'s are *P*'s.
- All cantaloupes are fruits.
- All fruits are seed-bearers.
- ∴ All seed-bearers are cantaloupes.

Intuitively, it is obvious that this is not a valid argument. The task of logic is to show why a syllogism is valid or not. An example of a valid syllogism is EIO in the second figure:

- No *P*'s are *M*'s.
- Some *S*'s are *M*'s.
- ∴ Some *S*'s are not *P*'s.
- No scientists are children.
- Some infants are children.
- ∴ Some infants are not scientists.

The validity of a syllogism depends on the relations among the classes referred to by the terms of the argument. If all of one class is contained in a second class and none of the second class is in a third, then none of the first class is in the third either. Using this principle and others like it, logicians have been able to establish which syllogisms are valid and which are not.

Arguments presented in ordinary contexts, even when statable in categorical propositions, may not be simple syllogisms. Often essential premises are not stated, because they are so obvious and trivial as not to require mentioning. When an essential premise is not stated, the argument is called an enthymeme. Enthymematic arguments need to have their hidden premises made explicit before a test for validity can be made. In addition, arguments often contain more than two premises. Indeed, some arguments can be structured as a sequence of syllogisms, where preliminary conclusions are expressly drawn and then are used as premises in later syllogisms. Such a chain of subarguments is called a sorites. The English logician and novelist Lewis Carroll devised clever, whimsical sorites that have entertained students for more than 100 years. For instance, in *Symbolic Logic* (1896) he presented the following argument, whose conclusion was left unexpressed:

All my sons are slim.
No child of mine is healthy who takes no exercise.
All gluttons who are children of mine are fat.
No daughter of mine takes any exercise.

In addition, certain crucial premises of this argument—such as “No slim persons are fat persons”—have not been expressed.

Other argument forms. The argument form most discussed and studied from the time of Aristotle to the early 19th century was the syllogism. But Aristotle himself noted that some arguments were expressed in propositions other than categorical ones. The following argument, for instance, has for its first premise a hypothetical proposition:

If all men are born equals, then all slaves are unjustly treated persons.
All men are born equals.
∴ All slaves are unjustly treated persons.

This is a hypothetical argument, often called a hypothetical syllogism. Hypothetical propositions have the form “If . . . , then —,” where the word “then” is often omitted. When, as above, the conclusion is obtained by the second premise’s affirming the antecedent, the argument is said to be by *modus ponens*. The conclusion in this case is the consequent of the hypothetical first premise.

A hypothetical argument can also be conducted by denying the consequent of the hypothetical premise and thereby concluding with a denial of the antecedent of the hypothetical. This form of hypothetical argument is called *modus tollens*, and the denials in either case are frequently expressed by the contradictory of the proposition at issue, either the antecedent or consequent of the hypothetical. An example of a *modus tollens* hypothetical argument is

If some persons are persons with rights to freedom, then all persons are persons with rights to freedom.
Not all persons are persons with rights to freedom.
∴ No persons are persons with rights to freedom.

Disjunctions are propositions in which the predicate is asserted to belong to one or another subject, or one or another predicate is asserted to belong to a subject: “Either *A*’s or *B*’s are *C*’s, or *A*’s are either *B*’s or *C*’s.” Another more complex disjunction takes two categorical propositions as alternatives: “Either *A*’s are *B*’s, or *C*’s are *D*’s.” A disjunctive argument (sometimes called a disjunctive syllogism) contains one of the three above disjunctive forms as one premise and the denial of one of the alternatives (disjuncts) as the second premise. The valid conclusion in these cases is the other alternative. A simple and traditional example is

Either God is unjust, or no men are eternally punished creatures.
God is not unjust.
∴ No men are eternally punished creatures.

The singular proposition here (“God is unjust”) is treated as a universal categorical proposition.

Sometimes the alternatives are meant to be exclusive—that is, if one is true, the other is false. When such is the case, a valid disjunctive argument can then be con-

structed by affirming one of the alternatives in a premise and subsequently concluding a denial of the other alternative. Thus,

Either Bacon or Shakespeare is the author of Hamlet.
Shakespeare is the author of *Hamlet*.
∴ Bacon is not the author of *Hamlet*.

Unfortunately, it is not always evident whether the disjunction is to be taken in the inclusive or the exclusive sense, and the careful logician will usually explicitly assert “*A* or *B*, but not both.” Examples of ambiguity of disjunction abound: “Newton or Leibniz is the discoverer of the calculus (possible codiscoverers);” “All diplomats are liars or failures.”

A combination of a disjunction and hypothetical propositions as premises gives rise to a type of argument known as a dilemma. The hypothetical propositions offer alternatives, either one of which leads to a (frequently unpalatable) conclusion. When the conclusions of both alternatives are the same, it is a simple dilemma; when they differ, it is a complex dilemma. If the antecedent of the hypothetical proposition is affirmed, and thus the consequent is also affirmed as conclusion, the argument is constructive. When the consequent is denied, and thus the antecedent is denied as conclusion, the argument is called destructive. Some illustrations of these types of dilemmas are displayed below. (For ease of reading, these propositions are not written in categorical form but are expressed as they would be colloquially.)

Simple constructive:

If a science furnishes useful facts, it is worthy of being cultivated; and if the study of it exercises the reasoning powers, it is worthy of being cultivated. But either a science furnishes useful facts, or its study exercises the reasoning powers. Therefore it is worthy of being cultivated.

(William Stanley Jevons, *Elementary Lessons in Logic* [1870].)

Complex constructive:

If there is censorship of the press, abuses of power will be concealed; and if there is no censorship, truth will be sacrificed to sensation. But there must either be censorship or not. Therefore either abuses of power will be concealed, or truth will be sacrificed to sensation.

(Horace William Brindley Joseph, *An Introduction to Logic* [1916].)

Destructive:

If this person were wise, he would not speak irreverently of Scripture in jest; and if he were good, he would not do so in earnest. But he does it either in jest or earnest. Therefore he is either not wise or not good.

(Richard Whately, *Elements of Logic* [1826].)

Symbolic logic. A number of developments during the Renaissance and immediately thereafter—the period of the emergence of modern science—led to increasing dissatisfaction with the traditional logic of the syllogism. In particular, the development of functional relations in natural science, the shift of interest from geometry to algebra in mathematics, the concern for the logical foundations of mathematics, and the call for a language that would reveal logical relations by its very notation (compare Gottfried Wilhelm Leibniz’ *characteristica universalis*) led to the developments in the 19th century that can be called the algebra of logic. It is notorious that the British mathematician and logician Augustus De Morgan (1847) found fault with the syllogism by pointing out that it cannot (easily) deal with the simple relational inference:

All horses are animals.
∴ All heads of horses are heads of animals.

Although various abbreviations were accomplished through symbols, even in the works of Aristotle himself, the use of symbols in an explicit formal system, the precursor of modern symbolic logic, began with George Boole (1847) and Ernst Schröder (1890–1905), was developed further by Gottlob Frege (1879), and finally culminated in the *Principia Mathematica* of Bertrand Russell and Alfred North Whitehead (1910–13). The formal systems of modern symbolic logic differ from earlier logical studies

Dilemmas

Hypothetical arguments

Modern systems of symbolic logic

that used symbols in that, in the former, totally artificial languages are rigorously developed using special symbols for precisely defined logical concepts. The rules of this language, both the syntactic rules for deduction and the semantic rules for interpreting expressions, are explicitly and precisely stated. The development of these symbolic formal systems within which deductive arguments can be represented yields a number of distinct advantages. A high degree of rigour can be attained. The sharp separation of semantics from syntax leads to a clear distinction between the validity of an argument (semantics) and the deducibility of the conclusion from axioms and premises (syntax). Additionally, the formal system, once made totally explicit, can itself be the object of study (see below *Logic systems: Metalogic*).

The logical relations among whole sentences is the basis of the modern symbolic approach. In effect, hypothetical and disjunctive arguments rather than the categorical syllogism become the centre of attention. Beginning with simple sentences that have no simpler sentences as components, one constructs compound sentences using sentential connectives. The truth value (either true or false) of the compound sentence depends then on the truth values of its components in a clear and explicit manner according to which function is represented by the sentential connective. For instance, the propositional truth function called conjunction, which is frequently represented by “ \cdot ” or “ $\&$,” has the value true when both the conjoined propositions have the value true; otherwise it has the value false. In other words, if p and q are arbitrary propositions, the sentence “ $p \cdot q$ ” represents a true proposition just in case both p and q are true propositions themselves. The formalization of these truth functions and the statement of the rules for inferring new sentences from earlier ones (the rules of inference) results in a formal system called the propositional calculus (PC).

Yet PC cannot deal with arguments formerly handled by the categorical syllogism. Some way of dealing with the internal structure of simple sentences needs to be developed. The great power of modern logic is based on the important notion of a propositional function. A propositional function acts on a domain of individuals and has the value true or false, depending on which individual (or individuals) is the argument of the function. Thus, “ $_$ is an even number” represents a propositional function whose value is true whenever the blank is filled by a numeral referring to an even number and false when the number is odd.

Instead of using expressions with blank spaces, which can be confusing if there is more than one blank, logicians utilize what are termed individual variables, expressions that hold open a place in a sentence fragment for the name of some individual. Individual variables are frequently lowercase letters from the end of the alphabet. So the example in the previous paragraph would be written: “ x is an even number.” This expression can become a sentence when the variable “ x ” is replaced by the name of some thing—a true sentence when that thing is an even number. There are other ways to convert such expressions into sentences. One can prefix the expression with a universal quantifier, “For all x .” Now the resulting sentence, “For all x , x is an even number,” expresses the false proposition that everything is an even number. Furthermore, prefixing the expression with an existential quantifier, “There is at least one x ,” yields the true sentence, “There is at least one thing such that it is an even number.”

Being an even number is a property that some individuals can have. Expressions that attribute a property to an individual are (monadic) predicates. It is customary to express simple predicates by uppercase letters placed before the individual term. Thus if E is used for the predicate “is an even number,” the expression Ex is intended to represent “ x is an even number.” Using monadic predicates, quantifiers, individual variables, and the sentential connectives developed in PC, it is possible to express all the categorical syllogisms and subsequently determine their validity. When rules of inference and possibly axioms are introduced, this system is called the monadic predicate calculus. When relations are asserted to hold between two

or more individuals, additional, n -adic, predicates enter the language. For example, using the uppercase letter L to express the dyadic relation of being less than, and taking a and b to be any (not necessarily different) numbers, one can assert that a is less than b by writing: Lab . The notation of dyadic relation symbols allows a simple expression, and solution, of De Morgan’s problem, mentioned above, about heads of horses. One may even introduce the notion of predicate variables; but, as long as there is no quantification over predicate variables, the resulting formal system is called the lower predicate calculus (LPC).

One further extension of LPC is usually made in modern logic. One special dyadic relation, represented by the equality sign, “ $=$,” placed between two terms, is taken to be the identity relation. Depending on the type of formal system that is being considered, either axioms of identity (e.g., “Everything is self-identical”) are adopted or else rules of inference governing transformations (e.g., “From any conclusion ϕ containing the name a and an earlier line of derivation, $a = b$, infer a new conclusion ϕ' containing b for some occurrences of a ”) are added to the earlier rules of the system. The resulting system, which in effect restricts the possible interpretations of LPC to the identity relation for the dyadic predicate “ $=$,” is called LPC with identity (or sometimes first-order logic with identity). Several considerations suggest that this is the most comprehensive logical system possible and that any other additions will no longer result in all logical truths, and only logical truths, as theorems.

In formal systems the emphasis shifts from arguments to deducing conclusions. The rules of inference of the system allow various transformations on, or inferences from, initial sequences of symbols. When no additional material assumptions are used, the final line of any such derivation is called a theorem of logic. When, however, assumptions about some field of inquiry are incorporated into the formal system, the theorems derived by using the rules of the system are theorems of the material theory. Thus, if certain postulates about the behaviour of moving bodies are laid down, one would derive theorems of kinematics—and similarly for arithmetic, geometry, and so on.

Modern logic in the last part of the 20th century can be divided into four major areas of investigation. The first area is proof theory, the study of the properties of formal systems and the derivations that can be accomplished within them. The second area is model theory, which investigates the various structures about which formal theories can be constructed. Here the emphasis is on what cannot be validly deduced from a set of material hypotheses. One attempts to find structures about which the hypotheses are true and yet for which a particular statement is false. Third is recursion theory, which deals with questions involving the decidability of the question of whether or not a sentence is deducible from a set of premises. This study has led to theories of computability, or the existence of mechanical procedures for solving problems associated with deducibility. Finally, there is the broad area of the foundations of mathematics, especially the logical grounding of the basic notions of set theory.

Applications of the formal methods of logic have burgeoned with the development of novel semantic devices such as “possible worlds.” It is now possible to provide a semantics for various modal logics dealing with such topics as necessarily true propositions, known propositions (as distinct from those merely believed), obligatory actions, and the structure of temporal relations. Previously, formulas of modal logic were merely uninterpreted sequences of symbols with no clear meanings. In addition, grammatical studies within the general field of linguistics has benefited from the seminal work of the American logician Richard Montague (1970) and subsequent developments.

Inductive logic. Inductive arguments intend to support their conclusions only to some degree; the premises do not necessitate the conclusion. Traditionally, the study of inductive logic was confined to either arguments by analogy or else methods of arriving at generalizations on the basis of a finite number of observations. A typical argument by analogy proceeds from the premise that two objects are observed to be similar with respect to a number of

Four major areas of investigation

Individual variables

attributes to the conclusion that the two objects are also similar with respect to another attribute. The strength of such arguments depends on the degree to which the attributes in question are related to each other.

The methods appropriate to inductive generalizations have been studied by modern philosophers from Francis Bacon in the early 17th century to William Whewell and John Stuart Mill in the 19th century. Proper inductive generalizations require that the observed instances referred to in the premises be obtained according to a careful method of varying the circumstances of observations, a rigorous search for exceptional cases, and attempts to

detect correlations or dependencies among the various phenomena.

In the 20th century, most notably in the work of Hans Reichenbach (1938), a distinction has been made between the context of discovery and the context of justification, between the nonlogical process for arriving at a general hypothesis and the logical relations that obtain between the hypothesis and the evidence for it—the so-called hypothetico-deductive method. In modern inductive logic, the probability calculus, or some variant of it, is called upon to explicate the notion of how observed evidence logically supports a theoretical hypothesis. (M.L.Sc.)

LOGIC SYSTEMS

Formal logic

The discipline known as formal logic takes as its main subject matter propositions (or statements, or assertively used sentences) and deductive arguments, and it abstracts from their content the structures or logical forms that they embody. The logician customarily uses a symbolic notation to express these structures clearly and unambiguously and to enable manipulations and tests of validity to be more easily applied. Although this discussion freely employs the technical notation of modern symbolic logic, its symbols are introduced gradually and with accompanying explanations so that the serious and attentive general reader should be able to follow the development of ideas.

Formal logic is an a priori, and not an empirical, study. In this respect it contrasts with the natural sciences and with all other disciplines that depend on observation for their data. Its nearest analogy is with pure mathematics; indeed, many logicians and pure mathematicians would regard their respective subjects as indistinguishable, or as merely two stages of the same unified discipline. Formal logic, therefore, is not to be confused with the empirical study of the processes of reasoning, which belongs to psychology. It must also be distinguished from the art of correct reasoning, which is the practical skill of applying logical principles to particular cases; and, even more sharply, it must be distinguished from the art of persuasion, in which invalid arguments are sometimes more effective than valid ones.

GENERAL OBSERVATIONS

Probably the most natural approach to formal logic is through the idea of the validity of an argument of the kind known as deductive. A deductive argument can be roughly characterized as one in which the claim is made that some proposition (the conclusion) follows with strict necessity from some other proposition or propositions (the premises)—*i.e.*, that it would be inconsistent or self-contradictory to assert the premises but deny the conclusion.

If a deductive argument is to succeed in establishing the truth of its conclusion, two quite distinct conditions must be met: first, the conclusion must really follow from the premises—*i.e.*, the deduction of the conclusion from the premises must be logically correct—and, second, the premises themselves must be true. An argument meeting both these conditions is called sound. Of these two conditions, the logician as such is concerned only with the first; the second, the determination of the truth or falsity of the premises, is the task of some special discipline or of common observation appropriate to the subject matter of the argument. When the conclusion of an argument is correctly deducible from its premises, the inference from the premises to the conclusion is said to be (deductively) valid, irrespective of whether the premises are true or false. Other ways of expressing the fact that an inference is deductively valid are to say that the truth of the premises gives (or would give) an absolute guarantee of the truth of the conclusion or that it would involve a logical inconsistency (as distinct from a mere mistake of fact) to suppose that the premises were true but the conclusion false.

The deductive inferences with which formal logic is concerned are, as the name suggests, those for which validity

depends not on any features of their subject matter but on their form or structure. Thus the two inferences

- Every dog is a mammal. (1)
Some quadrupeds are dogs.
∴ Some quadrupeds are mammals.

and

- Every anarchist is a believer in free love. (2)
Some members of the government party are anarchists.
∴ Some members of the government party are believers in free love.

differ in subject matter and hence require different procedures to check the truth or falsity of their premises. But their validity is ensured by what they have in common—namely, that the argument in each is of the form:

- Every X is a Y . (3)
Some Z 's are X 's.
∴ Some Z 's are Y 's.

Line (3) above may be called an inference form, and (1) and (2) are then instances of that inference form. The letters— X , Y , and Z —in (3) mark the places into which expressions of a certain type may be inserted. Symbols used for this purpose are known as variables; their use is analogous to that of the x in algebra, which marks the place into which a numeral can be inserted. An instance of an inference form is produced by replacing all the variables in it by appropriate expressions—*i.e.*, ones that make sense in the context—and by doing so uniformly—*i.e.*, by substituting the same expression wherever the same variable recurs. The feature of (3) that guarantees that every instance of it will be valid is its construction in such a manner that every uniform way of replacing its variables to make the premises true automatically makes the conclusion true also, or, in other words, that no instance of it can have true premises but a false conclusion. In virtue of this feature, the form (3) is termed a valid inference form. In contrast,

- Every X is a Y . (4)
Some Z 's are Y 's.
∴ Some Z 's are X 's.

is not a valid inference form, for, although instances of it can be produced in which premises and conclusion are all true, instances of it can also be produced in which the premises are true but the conclusion is false—*e.g.*,

- Every dog is a mammal. (5)
Some winged creatures are mammals.
∴ Some winged creatures are dogs.

Formal logic as a study is concerned with inference forms rather than with particular instances of them; one of its tasks is to discriminate between valid and invalid inference forms and to explore and systematize the relations that hold among valid ones.

Closely related to the idea of a valid inference form is that of a valid proposition form. A proposition form is an expression of which the instances (produced as before by appropriate and uniform replacements for variables) are not inferences from several propositions to a conclusion but rather propositions taken individually, and a valid

Valid proposition forms

A priori nature of formal logic

Inference forms

proposition form is one for which all of the instances are true propositions. A simple example is

$$\text{Nothing is both an } X \text{ and a non-}X. \quad (6)$$

Formal logic is concerned with proposition forms as well as with inference forms. The study of proposition forms can, in fact, be made to include that of inference forms in the following way: let the premises of any given inference form (taken together) be abbreviated by alpha (α) and its conclusion by beta (β). Then the condition stated above for the validity of the inference form " α , therefore β " amounts to saying that no instance of the proposition form " α and not- β " is true—i.e., that every instance of the proposition form

$$\text{Not both: } \alpha \text{ and not-}\beta \quad (7)$$

is true—or that line (7), fully spelled out, of course—is a valid proposition form. The study of proposition forms, however, cannot be similarly accommodated under the study of inference forms, and so for reasons of comprehensiveness it is usual to regard formal logic as the study of proposition forms. Because a logician's handling of proposition forms is in many ways analogous to a mathematician's handling of numerical formulas, the systems he constructs are often called calculi.

Much of the work of a logician proceeds at a more abstract level than that of the foregoing discussion. Even a formula such as (3) above, though not referring to any specific subject matter, contains expressions like "every" and "is a," which are thought of as having a definite meaning, and the variables are intended to mark the places for expressions of one particular kind (roughly, common nouns or class names). It is possible, however—and for some purposes it is essential—to study formulas without attaching even this degree of meaningfulness to them. The construction of a system of logic, in fact, involves two distinguishable processes: one consists in setting up a symbolic apparatus—a set of symbols, rules for stringing these together into formulas, and rules for manipulating these formulas; the second consists in attaching certain meanings to these symbols and formulas. If only the former is done, the system is said to be uninterpreted, or purely formal; if the latter is done as well, the system is said to be interpreted. This distinction is important, because systems of logic turn out to have certain properties quite independently of any interpretations that may be placed upon them. An axiomatic system of logic can be taken as an example—i.e., a system in which certain unproved formulas, known as axioms, are taken as starting points, and further formulas (theorems) are proved on the strength of these. As will appear later (*Axiomatization of PC*), the question whether a sequence of formulas in an axiomatic system is a proof or not depends solely on which formulas are taken as axioms and on what the rules are for deriving theorems from axioms, and not at all on what the theorems or axioms mean. Moreover, a given uninterpreted system is in general capable of being interpreted equally well in a number of different ways; hence, in studying an uninterpreted system, one is studying the structure that is common to a variety of interpreted systems. Normally a logician who constructs a purely formal system does have a particular interpretation in mind, and his motive for constructing it is the belief that when this interpretation is given to it the formulas of the system will be able to express true principles in some field of thought; but, for the above reasons among others, he will usually take care to describe the formulas and state the rules of the system without reference to interpretation and to indicate as a separate matter the interpretation that he has in mind.

Many of the ideas used in the exposition of formal logic, including some that are mentioned above, raise problems that belong to philosophy rather than to logic itself. Examples are: What is the correct analysis of the notion of truth? What is a proposition, and how is it related to the sentence by which it is expressed? Are there some kinds of sound reasoning that are neither deductive nor inductive? Fortunately, it is possible to learn to do formal logic without having satisfactory answers to such questions, just as it is possible to do mathematics without answering questions

belonging to the philosophy of mathematics such as: Are numbers real objects or mental constructs?

THE PROPOSITIONAL CALCULUS

Basic features of PC. The simplest and most basic branch of logic is the propositional calculus, hereafter called PC, named from the fact that it deals only with complete, unanalyzed propositions and certain combinations into which they enter. Various notations for PC are used in the literature. In that used here the symbols employed in PC first comprise variables (for which the letters p, q, r, \dots are used, with or without numerical subscripts); second, operators (for which the symbols " \sim ," " \cdot ," " \vee ," " \supset ," " \equiv " are employed); and third, brackets or parentheses. The rules for constructing formulas are discussed below (see *Formation rules for PC*), but the intended interpretations of these symbols—i.e., the meanings to be given to them—are indicated here immediately: The variables are to be viewed as representing unspecified propositions or as marking the places in formulas into which sentences, and only sentences, may be inserted. (This is sometimes expressed by saying that variables range over propositions, or that they take propositions as their values.) Hence they are often called propositional variables. It is assumed that every proposition is either true or false and that no proposition is both true and false. Truth and falsity are said to be the truth values of propositions. The function of an operator is to form a new proposition from one or more given propositions, called the arguments of the operator. The operators $\sim, \cdot, \vee, \supset,$ and \equiv correspond respectively to the English expressions "not," "and," "or," "if . . . , then" (or "implies"), and "is equivalent to," when these are used in the following senses:

1. Given a proposition p , then $\sim p$ ("not p ") is to count as false when p is true and true when p is false; " \sim " (when thus interpreted) is known as the negation sign, and $\sim p$ as the negation of p .
2. Given any two propositions p and q , then $p \cdot q$ (" p and q ") is to count as true when p and q are both true and as false in all other cases (namely, when p is true and q false, when p is false and q true, and when p and q are both false); $p \cdot q$ is said to be the conjunction of p and q ; " \cdot " is known as the conjunction sign, and its arguments (p, q) as conjuncts.
3. $p \vee q$ (" p or q ") is to count as false when p and q are both false and true in all other cases; thus it represents the assertion that at least one of p and q is true. $p \vee q$ is known as the disjunction of p and q ; " \vee " is the disjunction sign, and its arguments (p, q) are known as disjuncts.
4. $p \supset q$ ("if p [then] q " or " p [materially] implies q ") is to count as false when p is true and q is false and as true in all other cases; hence it has the same meaning as "either not- p or q " or as "not both; p and not- q ." The symbol " \supset " is known as the (material) implication sign, the first argument as the antecedent, and the second as the consequent; $q \supset p$ is known as the converse of $p \supset q$.
5. Finally, $p \equiv q$ (" p is [materially] equivalent to q " or " p if and only if q ") is to count as true when p and q have the same truth value (i.e., either when both are true or when both are false), and false when they have different truth values; the arguments of " \equiv " (the [material] equivalence sign) are called equivalents.

Brackets are used to indicate grouping; they make it possible to distinguish, for example, between $p \cdot (q \vee r)$ ("both p and either- q -or- r ") and $(p \cdot q) \vee r$ ("either both- p -and- q or r "). Precise rules for bracketing are given below.

All PC operators take propositions as their arguments, and the result of applying them is also in each case a proposition. For this reason they are sometimes called proposition-forming operators on propositions or, more briefly, propositional connectives. An operator that, like \sim , requires only a single argument is known as a monadic operator; operators that, like all the others listed, require two arguments are known as dyadic.

All PC operators also have the following important characteristic: given the truth values of the arguments, the truth value of the proposition formed by them and the

Interpreted and uninterpreted systems

Fundamental definitions

Truth function-ality

operator is determined in every case. An operator that has this characteristic is known as a truth-functional operator, and a proposition formed by such an operator is called a truth function of the operator's argument(s). The truth functionality of the PC operators is clearly brought out by summarizing the above account of them in Table 1. In it, "true" is abbreviated by "1" and "false" by "0," and to the left of the vertical line are tabulated all possible combinations of truth values of the operators' arguments. The columns of 1s and 0s under the various truth functions indicate their truth values for each of the cases; these columns are known as the truth tables of the relevant operators. It should be noted that any column of four 1s or 0s or both will specify a dyadic truth-functional operator. Because there are precisely 2^4 (i.e., 16) ways of forming a string of four symbols each of which is to be either 1 or 0 (1111, 1110, 1101, . . . 0000), there are 16 such operators in all; the four that are listed here are only the four most generally useful ones.

Table 1: Truth Table for Most Common Operators

monadic operator		dyadic operators				
p	$\sim p$	$p \cdot q$	$p \vee q$	$p \supset q$	$p \equiv q$	
1	0	1 1	1	1	1	1
0	1	1 0	0	1	0	0
		0 1	0	1	1	0
		0 0	0	0	1	1

Formation rules for PC. In any system of logic it is necessary to specify which sequences of symbols are to count as acceptable formulas—or, as they are usually called, well-formed formulas (wffs). Rules that specify this are called formation rules. From an intuitive point of view, it is desirable that the wffs of PC be just those sequences of PC symbols that, in terms of the interpretation given above, make sense and are unambiguous; and this can be ensured by stipulating that the wffs of PC are to be all those expressions constructed in accordance with the following PC formation rules, and only these:

- FR1. A variable standing alone is a wff.
- FR2. If α is a wff, so is $\sim\alpha$.
- FR3. If α and β are wffs, $(\alpha \cdot \beta)$, $(\alpha \vee \beta)$, $(\alpha \supset \beta)$, and $(\alpha \equiv \beta)$ are wffs.

(In these rules α and β are variables representing arbitrary formulas of PC. They are not themselves symbols of PC but are used in discoursing about PC. Such variables are known as metalogical variables. For further explanation, see below *Metalogic*.) It should be noted that the rules, though designed to ensure unambiguous sense for the wffs of PC under the intended interpretation, are themselves stated without any reference to interpretation and in such a way that there is an effective procedure for determining, again without any reference to interpretation, whether any arbitrary string of symbols is a wff or not. (An effective procedure is one that is "mechanical" in nature and can always be relied on to give a definite result in a finite number of steps. The notion of effectiveness plays an important role in formal logic.)

Examples of wffs are: p ; $\sim q$; $\sim(p \cdot q)$, "not both p and q "; $[\sim p \vee (q \equiv p)]$, "either not p or else q is equivalent to p ."

For greater ease in writing or reading formulas, the formation rules are often relaxed. The following relaxations are common: (1) Brackets enclosing a complete formula may be omitted. (2) The typographical style of brackets may be varied within a formula to make the pairing of brackets more evident to the eye. (3) Conjunctions and disjunctions may be allowed to have more than two arguments—for example, $p \cdot (q \supset r) \cdot \sim r$ may be written instead of $[p \cdot (q \supset r)] \cdot \sim r$. (The conjunction $p \cdot q \cdot r$ is then interpreted to mean that p , q , and r are all true, $p \vee q \vee r$ to mean that at least one of p , q , and r is true, and so forth.)

Validity in PC. Given the standard interpretation, a wff of PC becomes a sentence, true or false, when all its variables are replaced by actual sentences. Such a wff is therefore a proposition form in the sense explained above

and hence is valid if and only if all its instances express true propositions. A wff of which all instances are false is said to be unsatisfiable, and one with some true and some false instances is said to be contingent.

An important problem for any logical system is the decision problem for the class of valid wffs of that system (sometimes simply called the decision problem for the system). This is the problem of finding an effective procedure, in the sense explained above in *Formation rules for PC*, for testing the validity of any wff of the system. Such a procedure is called a decision procedure. For some systems a decision procedure can be found; the decision problem for a system of this sort is then said to be solvable, and the system is said to be a decidable one. For other systems it can be proved that no decision procedure is possible; the decision problem for such a system is then said to be unsolvable, and the system is said to be an undecidable one.

PC is a decidable system. In fact, several decision procedures for it are known. Of these the simplest and most important theoretically (though not always the easiest to apply in practice) is the method of truth tables, which will now be explained briefly. Since all the operators in a wff of PC are truth-functional, in order to discover the truth value of any instance of such a wff, it is unnecessary to consider anything but the truth values of the sentences replacing the variables. In other words, the assignment of a truth value to each of the variables in a wff uniquely determines a truth value for the whole wff. Since there are only two truth values and each wff contains only a finite number of variables, there are only a finite number of truth-value assignments to the variables to be considered (if there are n distinct variables in the wff, there are 2^n such assignments); these can easily be systematically tabulated. For each of these assignments the truth tables for the operators then enable one to calculate the resulting truth value of the whole wff; and if and only if this truth value is truth in each case is the wff valid. As an example, $[(p \supset q) \cdot r] \supset [(\sim r \vee p) \supset q]$ may be tested for validity. This formula states that "if one proposition implies a second one, and a certain third proposition is true, then if either that third proposition is false or the first is true, the second is true."

The calculation is shown in Table 2. As before, 1 represents truth and 0 falsity. Since the wff contains three variables, there are 2^3 (i.e., 8) different assignments to the variables to be considered, which therefore generate the eight lines of the table. These assignments are tabulated to the left of the vertical line. The numbers in parentheses at the foot indicate the order in which the steps (from 1 through 6) are to be taken in determining the truth values (1 or 0) to be entered in the table. Thus column 1, falling under the symbol \supset , sets down the values of $p \supset q$ for each assignment, obtained from the columns under p and q by the truth table for \supset ; column 2, for $(p \supset q) \cdot r$, is then obtained by employing the values in column 1 together with those in the column under r by use of the truth table for \cdot ; . . . until finally column 6, which gives the values for the whole wff, is obtained from columns 2 and 5. This column is called the truth table of the whole wff. Since it consists entirely of 1s, it shows that the wff is true for every assignment given to the variables and is therefore valid. A wff for which the truth table consists entirely of 0s is never satisfied, and a wff for which the truth table contains at least one 1 and at least one 0 is contingent. It follows from the formation rules and from the fact that an initial truth table has been specified for each operator that

Decision procedure with truth tables

Well-formed formulas (wffs)

Table 2: Test for Validity by Truth Table

p	q	r	$[(p \supset q) \cdot r] \supset [(\sim r \vee p) \supset q]$					
1	1	1	1	1	1	0	1	1
1	1	0	1	0	1	1	1	1
1	0	1	0	0	1	0	1	0
1	0	0	0	0	1	1	1	0
0	1	1	1	1	1	0	0	1
0	1	0	1	0	1	1	1	1
0	0	1	1	1	1	0	0	1
0	0	0	1	0	1	1	1	0
			(1)	(2)	(6)	(3)	(4)	(5)

a truth table can be constructed for any given wff of PC. Among the more important valid wffs of PC are those of Table 3, all of which can be shown to be valid by a mechanical application of the truth-table method. They can also be seen to express intuitively sound general principles about propositions. For instance, because “not (. . . or —)” can be rephrased as “neither . . . nor —,” the first De Morgan law can be read as “both p and q if and only if neither not- p nor not- q ”; thus it expresses the principle that two propositions are jointly true if and only if neither of them is false.

Construction of valid wffs

law	formula
Law of identity	$p \equiv p$
Law of double negation	$p \equiv \sim\sim p$
Law of excluded middle	$p \vee \sim p$
Law of noncontradiction	$\sim(p \cdot \sim p)$
De Morgan laws	$(p \cdot q) \equiv \sim(\sim p \vee \sim q)$ $(p \vee q) \equiv \sim(\sim p \cdot \sim q)$
Commutative laws	$(p \vee q) \equiv (q \vee p)$ $(p \cdot q) \equiv (q \cdot p)$
Associative laws	$[(p \vee q) \vee r] \equiv [p \vee (q \vee r)]$ $[(p \cdot q) \cdot r] \equiv [p \cdot (q \cdot r)]$
Law of transposition	$(p \supset q) \equiv (\sim q \supset \sim p)$
Distributive laws	$[p \cdot (q \vee r)] \equiv [(p \cdot q) \vee (p \cdot r)]$ $[p \vee (q \cdot r)] \equiv [(p \vee q) \cdot (p \vee r)]$
Law of permutation	$[p \supset (q \supset r)] \equiv [q \supset (p \supset r)]$
Law of syllogism	$(p \supset q) \supset [(q \supset r) \supset (p \supset r)]$
Law of importation	$[p \supset (q \supset r)] \supset [(p \cdot q) \supset r]$
Law of exportation	$[(p \cdot q) \supset r] \supset [p \supset (q \supset r)]$

Whenever, as is the case in most of the examples given, a wff of the form $\alpha \equiv \beta$ is valid, the corresponding wffs $\alpha \supset \beta$ and $\beta \supset \alpha$ are also valid. For instance, because $(p \cdot q) \equiv \sim(\sim p \vee \sim q)$ is valid, so are $(p \cdot q) \supset \sim(\sim p \vee \sim q)$ and $\sim(\sim p \vee \sim q) \supset (p \cdot q)$.

Moreover, although $p \supset q$ does not mean that q can be deduced from p , yet whenever a wff of the form $\alpha \supset \beta$ is valid, the inference form “ α , therefore β ” is likewise valid. This fact is easily seen from the fact that $\alpha \supset \beta$ means the same as “not both: α and not- β ”; for, as was noted above, whenever the latter is a valid proposition form, “ α , therefore β ” is a valid inference form.

Let α be any wff. If any variable in it is now uniformly replaced by some wff, the resulting wff is called a substitution-instance of α . Thus $[p \supset (q \vee \sim r)] \equiv [\sim(q \vee \sim r) \supset \sim p]$ is a substitution-instance of $(p \supset q) \equiv (\sim q \supset \sim p)$, obtained from it by replacing q uniformly by $(q \vee \sim r)$. It is an important principle that, whenever a wff is valid, so is every substitution-instance of it (the rule of [uniform] substitution).

A further important principle is the rule of substitution of equivalents. Two wffs, α and β , are said to be equivalents when $\alpha \equiv \beta$ is valid. (The wffs α and β are equivalents if and only if they have identical truth tables.) The rule states that, if any part of a wff is replaced by an equivalent of that part, the resulting wff and the original are also equivalents. Such replacements need not be uniform. The application of this rule is said to make an equivalence transformation.

Interdefinability of operators. The rules that have just been stated would enable the first De Morgan law listed in Table 3 to transform any wff containing any number of occurrences of \cdot into an equivalent wff in which \cdot does not appear at all but in place of it certain complexes of \sim and \vee arise. Similarly, since $\sim p \vee q$ has the same truth table as $p \supset q$, $(p \supset q) \equiv (\sim p \vee q)$ is valid, and any wff containing \supset can therefore be transformed into an equivalent wff containing \sim and \vee but not \supset . And, since $(p \equiv q) \equiv [(p \supset q) \cdot (q \supset p)]$ is valid, any wff containing \equiv can be transformed into an equivalent containing \supset and \cdot but not \equiv , and thus in turn by the previous steps it can be further transformed into one containing \sim and \vee but neither \equiv nor \supset nor \cdot . Thus, for every wff of PC there is an equivalent wff, expressing precisely the same truth function, in which the only operators are \sim and \vee , though the meaning of this wff will usually be much less clear than that of the original.

An alternative way of presenting PC, therefore, is to begin

with the operators \sim and \vee only and to define the others in terms of these. The operators \sim and \vee are then said to be primitive. If “ \equiv_{Df} ” is used to mean “is defined as,” then the relevant definitions can be set down as follows:

$$\begin{aligned} (\alpha \cdot \beta) &=_{\text{Df}} \sim(\sim\alpha \vee \sim\beta) \\ (\alpha \supset \beta) &=_{\text{Df}} (\sim\alpha \vee \beta) \\ (\alpha \equiv \beta) &=_{\text{Df}} [(\alpha \supset \beta) \cdot (\beta \supset \alpha)] \end{aligned}$$

in which α and β are any wffs of PC. These definitions are not themselves wffs of PC, nor is \equiv_{Df} a symbol of PC; they are metalogical statements about PC, used to introduce the new symbols \cdot , \supset , and \equiv into the system. If PC is regarded as a purely uninterpreted system, the expression on the left in a definition is simply a convenient abbreviation of the expression on the right. If, however, PC is thought of as having its standard interpretation, the meanings of \sim and \vee will first of all have been stipulated by truth tables, and then the definitions will lay it down that the expression on the left is to be understood as having the same meaning (*i.e.*, the same truth table) as the expression on the right. It is easy to check that the truth tables obtained in this way for \cdot , \supset , and \equiv are precisely the ones that were originally stipulated for them.

An alternative to taking \sim and \vee as primitive is to take \sim and \cdot as primitive and to define $(\alpha \vee \beta)$ as $\sim(\sim\alpha \cdot \sim\beta)$, to define $(\alpha \supset \beta)$ as $\sim(\alpha \cdot \sim\beta)$, and to define $(\alpha \equiv \beta)$ as before. Yet another possibility is to take \sim and \supset as primitive and to define $(\alpha \vee \beta)$ as $(\sim\alpha \supset \beta)$, $(\alpha \cdot \beta)$ as $\sim(\alpha \supset \sim\beta)$, and $(\alpha \equiv \beta)$ as before. In each case, precisely the same wffs that were valid in the original presentation of the system are still valid.

Axiomatization of PC. The basic idea of constructing an axiomatic system is that of choosing certain wffs (known as axioms) as starting points and giving rules for deriving further wffs (known as theorems) from them. Such rules are called transformation rules. Sometimes the word “theorem” is used to cover axioms as well as theorems; the word “thesis” is also used for this purpose.

An axiomatic basis consists of:

- (1) a list of primitive symbols, together with any definitions that may be thought convenient,
- (2) a set of formation rules, specifying which sequences of symbols are to count as wffs,
- (3) a list of wffs selected as axioms, and
- (4) a set of (one or more) transformation rules, which enable new wffs (theorems) to be obtained by performing certain specified operations on axioms or previously obtained theorems.

Axiomatic basis and interpretation

Definitions, where they occur, can function as additional transformation rules, to the effect that, if in any theorem any expression of the form occurring on one side of a definition is replaced by the corresponding expression of the form occurring on the other side, the result is also to count as a theorem. A proof or derivation of a wff α in an axiomatic system S is a sequence of wffs of which the last is α itself and each wff in the sequence is either an axiom of S or is derived from some axiom(s) or some already-derived theorem(s) or both by one of the transformation rules of S . A wff is a theorem of S if and only if there is a proof of it in S .

Care is usually taken, in setting out an axiomatic basis, to avoid all reference to interpretation. It must be possible to tell purely from the construction of a wff whether it is an axiom or not. Moreover, the transformation rules must be so formulated that there is an effective way of telling whether any purported application of them is a correct application or not, and hence whether a purported proof of a theorem really is a proof or not. An axiomatic system will then be a purely formal structure, on which any one of a number of interpretations, or none at all, may be imposed without affecting the question of which wffs are theorems. Normally, however, an axiomatic system is constructed with a certain interpretation in mind: the transformation rules are so formulated that under that interpretation they are validity-preserving (*i.e.*, the results of applying them to valid wffs are always themselves valid wffs); and the chosen axioms either are valid wffs or are expressions of principles of which it is desired to explore the consequences.

Properties of an axiomatic system: PM

Probably the best-known axiomatic system for PC is the following one, which, since it is derived from *Principia Mathematica*, by Whitehead and Russell, is often called PM:

Primitive symbols: $\sim, \vee, (\cdot),$ and an infinite set of variables, p, q, r, \dots (with or without numerical subscripts).

Definitions of \cdot, \supset, \equiv as above (*Interdefinability of operators*).

Formation rules as above (*Formation rules for PC*), except that formation rule 3 can be abbreviated to "If α and β are wffs, $(\alpha \vee \beta)$ is a wff," since $\cdot, \supset,$ and \equiv are not primitive.

Axioms:

1. $(p \vee p) \supset p$
2. $q \supset (p \vee q)$
3. $(p \vee q) \supset (q \vee p)$
4. $(q \supset r) \supset [(p \vee q) \supset (p \vee r)]$

Axiom 4 can be read, "If q implies r , then, if either p or q , either p or r ."

Transformation rules:

1. The result of uniformly replacing any variable in a theorem by any wff is a theorem (rule of substitution).
2. If α and $(\alpha \supset \beta)$ are theorems, then β is a theorem (rule of detachment, or modus ponens).

Relative to a given criterion of validity, an axiomatic system is sound if every theorem is valid, and it is complete (or, more specifically, weakly complete) if every valid wff is a theorem. The axiomatic system PM can be shown to be both sound and complete relative to the criterion of validity given in *Validity in PC* above.

An axiomatic system is consistent if, whenever a wff α is a theorem, $\sim\alpha$ is not a theorem. (In terms of the standard interpretation, this means that no pair of theorems can ever be derived one of which is the negation of the other.) It is strongly complete if the addition to it (as an extra axiom) of any wff whatever that is not already a theorem would make the system inconsistent. Finally, an axiom or transformation rule is independent (in a given axiomatic system) if it cannot be derived from the remainder of the basis (or—which comes to the same thing—if its omission from the basis would make the derivation of certain theorems impossible). It can, moreover, be shown that PM is consistent and strongly complete and that each of its axioms and transformation rules is independent.

A considerable number of other axiomatic bases for PC, each having all the above properties, are known. The task of proving that they have these properties belongs to metalogic.

In some standard expositions of formal logic, the place of axioms is taken by axiom schemata, which, instead of presenting some particular wff as an axiom, lay it down that any wff of a certain form is an axiom. For example, in place of axiom 1 in PM one might have the axiom schema "Every wff of the form $(\alpha \vee \alpha) \supset \alpha$ is an axiom"; and analogous schemata can be substituted for the other axioms. The number of axioms would then become infinite, but, on the other hand, the rule of substitution would no longer be needed, and modus ponens could be the only transformation rule. This method makes no difference to the theorems that can be derived; but in some branches of logic (though not in PC) it is simpler to work with axiom schemata rather than with particular axioms and substitution rules. Having an infinite number of axioms causes no trouble provided that there is an effective way of telling whether a wff is an axiom or not.

Special systems of PC. *Partial systems of PC.* Various propositional calculi have been devised to express a narrower range of truth functions than those of PC as expounded above. Of these, the one that has been most fully studied is the pure implicational calculus (PIC), in which the only operator is \supset , and the wffs are precisely those wffs of PC that can be built up from variables, \supset , and brackets alone. Formation rules 2 and 3 of *Formation rules for PC* (above) are therefore replaced by the rule that if α and β are wffs, $(\alpha \supset \beta)$ is a wff. As in ordinary PC, $p \supset q$ is interpreted as "p materially implies q"—i.e., as true except when p is true but q false. The truth-table

test of validity can then be straightforwardly applied to wffs of PIC.

The task of axiomatizing PIC is that of finding a set of valid wffs, preferably few in number and relatively simple in structure, from which all other valid wffs of the system can be derived by straightforward transformation rules. The best-known basis, which was formulated in 1930, has the transformation rules of substitution and modus ponens (as in PM) and the following axioms:

1. $p \supset (q \supset p)$
2. $[(p \supset q) \supset p] \supset p$
3. $(p \supset q) \supset [(q \supset r) \supset (p \supset r)]$

Axioms 1 and 3 are closely related to axioms 2 and 4 of PM respectively (see above *Axiomatization of PC*). It can be shown that the basis is complete and that each axiom is independent.

Under the standard interpretation, the above axioms can be thought of as expressing the following principles: (1) "If a proposition p is true, then if some arbitrary proposition q is true, p is (still) true." (2) "If the fact that a certain proposition p implies some arbitrary proposition q implies that p itself is true, then p is (indeed) true." (3) "If one proposition (p) implies a second (q), then if that second proposition implies a third (r), the first implies the third." The completeness of the basis is, however, a formal matter, not dependent on these or any other readings of the formulas.

An even more economical complete basis for PIC contains the same transformation rules but the sole axiom

$$[(p \supset q) \supset r] \supset [(r \supset p) \supset (s \supset p)].$$

It has been proved that this is the shortest possible single axiom that will give a complete basis for PIC with these transformation rules.

Since PIC contains no negation sign, the previous account of consistency is not significantly applicable to it. Alternative accounts of consistency have, however, been proposed, according to which a system is consistent (1) if no wff consisting of a single variable is a theorem or (2) if not every wff is a theorem. The bases stated are consistent in these senses.

Nonstandard versions of PC. Qualms have sometimes been expressed about the intuitive soundness of some formulas that are valid in "orthodox" PC, and these qualms have led some logicians to construct a number of propositional calculi that deviate in various ways from PC as expounded above.

Underlying ordinary PC is the intuitive idea that every proposition is either true or false, an idea that finds its formal expression in the stipulation that variables shall have two possible values only—namely, 1 and 0. (For this reason the system is often called the two-valued propositional calculus.) This idea has been challenged on various grounds. Following a suggestion made by Aristotle, some logicians have maintained that propositions about those events in the future that may or may not come to pass are neither true nor false but "neuter" in truth value. Aristotle's example, which has received much discussion, is "There will be a sea battle tomorrow." It has also been maintained, by the English philosopher Sir Peter Strawson and others, that, for propositions with subjects that do not have anything actual corresponding to them—such as "The present king of France is wise" (assuming that France has no king) or "All John's children are asleep" (assuming that John has no children)—the question of truth or falsity "does not arise." Another view is that a third truth value (say "half-truth") ought to be recognized as existing intermediate between truth and falsity; thus it has been advanced that certain familiar states of the weather make the proposition "It is raining" neither definitely true nor definitely false but something in between the two.

The issues raised by the above examples no doubt differ significantly, but they all suggest a threefold rather than a twofold division of propositions and hence the possibility of a logic in which the variables may take any of three values (say 1, $1/2$, and 0), with a consequent revision of the standard PC account of validity. Several such three-valued logics have been constructed and investigated; a brief account will be given here of one of them, in which the

Axiom schemata

The true-false dichotomy: its problems

most natural interpretation of the additional value ($\frac{1}{2}$) is as "half-true," with 1 and 0 representing truth and falsity as before. The formation rules are as they were for orthodox PC, but the meaning of the operators is extended to cover cases in which at least one argument has the value $\frac{1}{2}$ by the five entries in Table 4. (Adopting one of the three values of the first argument, p , given in the leftmost column [1, $\frac{1}{2}$, or 0] and, for the dyadic operators, one of the three values of the second, q , in the top row—above the line—one then finds the value of the whole formula by reading across for p and down for q .) It will be seen that these tables, due to the Polish logician Jan Łukasiewicz, are the same as the ordinary two-valued ones when the arguments have the values 1 and 0. The other values are intended to be intuitively plausible extensions of the principles underlying the two-valued calculus to cover the cases involving half-true arguments. Clearly, these tables enable a person to calculate a determinate value (1, $\frac{1}{2}$, or 0) for any wff, given the values assigned to the variables in it; a wff is valid in this calculus if it has the value 1 for every assignment to its variables. Since the values of formulas when the variables are assigned only the values 1 and 0 are the same as in ordinary PC, every wff that is valid in the present calculus is also valid in PC. Some wffs that are valid in PC are, however, now no longer valid. An example is $(p \vee \sim p)$, which, when p has the value $\frac{1}{2}$, also has the value $\frac{1}{2}$. This reflects the idea that if one admits the possibility of a proposition's being half-true, he can no longer hold of every proposition without restriction that either it or its negation is true.

Given the truth tables for the operators in Table 4, it is possible to take \sim and \supset as primitive and to define $(\alpha \vee \beta)$ as $[(\alpha \supset \beta) \supset \beta]$ —though not as $(\sim \alpha \supset \beta)$ as in ordinary PC; $(\alpha \cdot \beta)$ as $\sim(\sim \alpha \vee \sim \beta)$; and $(\alpha \equiv \beta)$ as $[(\alpha \supset \beta) \cdot (\beta \supset \alpha)]$. With these definitions as given, all valid wffs constructed from variables and $\sim, \cdot, \vee, \supset,$ and \equiv can be derived by substitution and modus ponens from the following four axioms:

1. $p \supset (q \supset p)$
2. $(p \supset q) \supset [(q \supset r) \supset (p \supset r)]$
3. $[(p \supset \sim p) \supset p] \supset p$
4. $(\sim p \supset \sim q) \supset (q \supset p)$

Other three-valued logics can easily be constructed. For example, the above tables might be modified so that $\sim \frac{1}{2} = \frac{1}{2}$, $\frac{1}{2} \supset 0 = \frac{1}{2}$, $\frac{1}{2} \equiv 0 = 0$, and $0 \equiv \frac{1}{2} = \frac{1}{2}$ all have the value 0 instead of $\frac{1}{2}$ as before, leaving everything else unchanged. The same definitions are then still possible, but the list of valid formulas is different; e.g., $\sim \sim p \supset p$, which was previously valid, now has the value $\frac{1}{2}$ when p has the value $\frac{1}{2}$. This system can also be successfully axiomatized. Other calculi with more than three values can also be constructed along analogous lines.

Other nonstandard calculi have been constructed by beginning with an axiomatization instead of a definition of validity. Of these, the best-known is the intuitionistic calculus, devised by Arend Heyting, one of the chief representatives of the intuitionist school of mathematicians, a group of theorists who deny the validity of certain types of proof used in classical mathematics (see MATHEMATICS, THE FOUNDATIONS OF). At least in certain contexts, members of this school regard the demonstration of the falsity of the negation of a proposition (a proof by reductio ad absurdum) as insufficient to establish the truth of the proposition in question. Thus they regard $\sim \sim p$ as an inadequate premise from which to deduce p and hence do

not accept the validity of the law of double negation in the form $\sim \sim p \supset p$. They do, however, regard a demonstration that p is true as showing that the negation of p is false and hence accept $p \supset \sim \sim p$ as valid. For somewhat similar reasons, these mathematicians also refuse to accept the validity of arguments based on the law of excluded middle $(p \vee \sim p)$. The intuitionistic calculus aims at presenting in axiomatic form those and only those principles of propositional logic that are accepted as sound in intuitionist mathematics. In this calculus, $\sim, \cdot, \vee,$ and \supset are all primitive, the transformation rules as before are substitution and modus ponens, and the axioms are the following:

1. $p \supset (p \cdot p)$
2. $(p \cdot q) \supset (q \cdot p)$
3. $(p \supset q) \supset [(p \cdot r) \supset (q \cdot r)]$
4. $[(p \supset q) \cdot (q \supset r)] \supset (p \supset r)$
5. $p \supset (q \supset p)$
6. $[p \cdot (p \supset q)] \supset q$
7. $p \supset (p \vee q)$
8. $(p \vee q) \supset (q \vee p)$
9. $[(p \supset r) \cdot (q \supset r)] \supset [(p \vee q) \supset r]$
10. $\sim p \supset (p \supset q)$
11. $[(p \supset q) \cdot (p \supset \sim q)] \supset \sim p$

From this basis neither $p \vee \sim p$ nor $\sim \sim p \supset p$ can be derived, though $p \supset \sim \sim p$ can. In this respect this calculus resembles the second of the three-valued logics described above. It is, however, not possible to give a truth-table account of validity—no matter how many values are used—that will bring out as valid precisely those wffs that are theorems of the intuitionistic calculus and no others.

Natural deduction method in PC. PC is often presented by what is known as the method of natural deduction. Essentially this consists of a set of rules for drawing conclusions from hypotheses (assumptions, premises) represented by wffs of PC and thus for constructing valid inference forms. It also provides a method of deriving from these inference forms valid proposition forms, and in this way it is analogous to the derivation of theorems in an axiomatic system. One such set of rules is presented in Table 5 (and there are various other sets that yield the same results).

Rules of inference and their application

Table 5: Sample Set of Rules for the Natural Deduction Method in Propositional Calculus

rule	given	one may then conclude
1. Modus ponens	α and $\alpha \supset \beta$	β
2. Modus tollens	$\sim \beta$ and $\alpha \supset \beta$	$\sim \alpha$
3. Double negation	α	$\sim \sim \alpha$
	$\sim \sim \alpha$	α
4. Conjunction introduction	α and β	$\alpha \cdot \beta$
5. Conjunction elimination	$\alpha \cdot \beta$	α and also β
6. Disjunction introduction	either α or β separately	$\alpha \vee \beta$
7. Disjunction elimination	$\alpha \vee \beta$, a derivation of γ from α , and a derivation of γ from β	γ
8. Conditional proof	a derivation of β from the hypothesis α (perhaps with the help of other hypotheses)	$\alpha \supset \beta$ as a conclusion from these other hypotheses (if any)
9. Reductio ad absurdum	a derivation of $\beta \cdot \sim \beta$ from the hypothesis α (perhaps with the help of other hypotheses)	$\sim \alpha$ as a conclusion from these other hypotheses (if any)

A natural deduction proof is a sequence of wffs beginning with one or more wffs as hypotheses; fresh hypotheses may also be added at any point in the course of a proof. The rules may be applied to any wff or group of wffs, as appropriate, that have already occurred in the sequence. In the case of rules 1–7, the conclusion is said to depend on all of those hypotheses that have been used in the series of applications of the rules that have led to this conclusion; i.e., it is claimed simply that the conclusion follows from these hypotheses, not that it holds in its own right. An application of rule 8 or rule 9, however, reduces by one the number of hypotheses on which the conclusion depends; and a hypothesis so eliminated is said to be a discharged hypothesis. In this way a wff may be reached that depends on no hypotheses at all. Such a wff is a theorem of logic. It can be shown that those theorems derivable by the rules stated above—together with the definition of $\alpha \equiv \beta$

Validity of wffs in three-valued logic

Intuitionistic calculus

Table 4: Truth Values for Common Operators in a Three-Valued Logic

p	q	negation ($\sim p$)	conjunction ($p \cdot q$)	disjunction ($p \vee q$)	implication ($p \supset q$)	equivalence ($p \equiv q$)
1	1	0	1	1	1	1
1	$\frac{1}{2}$	0	$\frac{1}{2}$	1	$\frac{1}{2}$	$\frac{1}{2}$
1	0	0	0	1	0	0
$\frac{1}{2}$	1	$\frac{1}{2}$	$\frac{1}{2}$	1	$\frac{1}{2}$	$\frac{1}{2}$
$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	1	1
$\frac{1}{2}$	0	$\frac{1}{2}$	0	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$
0	1	1	0	1	1	0
0	$\frac{1}{2}$	1	0	$\frac{1}{2}$	1	$\frac{1}{2}$
0	0	1	0	0	1	1

as $(\alpha \supset \beta) \cdot (\beta \supset \alpha)$ —are precisely the valid wffs of PC. A set of natural deduction rules yielding as theorems all the valid wffs of a system is complete (with respect to that system) in a sense obviously analogous to that in which an axiomatic basis was said to be complete in *Axiomatization of PC* (above).

Sample proof of a theorem of logic

As an illustration, the formula $[(p \supset q) \cdot (p \supset r)] \supset [p \supset (q \cdot r)]$ will be derived as a theorem of logic by the natural deduction method. (The sense of this formula is that, if a proposition $[p]$ implies each of two other propositions $[q, r]$, then it implies their conjunction.) Explanatory comments follow the proof.

1	(1)	$(p \supset q) \cdot (p \supset r)$	hypothesis
2	(2)	p	hypothesis
1	(3)	$p \supset q$	1, conjunction elimination
1	(4)	$p \supset r$	1, conjunction elimination
1, 2	(5)	q	2, 3, modus ponens
1, 2	(6)	r	2, 4, modus ponens
1, 2	(7)	$q \cdot r$	5, 6, conjunction introduction
1	(8)	$p \supset (q \cdot r)$	2, 7, conditional proof
	(9)	$[(p \supset q) \cdot (p \supset r)] \supset [p \supset (q \cdot r)]$	1, 8, conditional proof

The figures in parentheses immediately preceding the wffs are simply for reference. To the right is indicated either that the wff is a hypothesis or that it is derived from the wffs indicated by the rules stated. On the left are noted the hypotheses on which the wff in question depends (either the first or the second line of the derivation, or both). Note that since 8 is derived by conditional proof from hypothesis 2 and from 7, which is itself derived from hypotheses 1 and 2, 8 depends only on hypothesis 1, and hypothesis 2 is discharged. Similarly, 9 depends on no hypotheses and is therefore a theorem.

By varying the above rules it is possible to obtain natural deduction systems corresponding to other versions of PC. For example, if the second part of the double negation rule is omitted and the rule is added that, given $\alpha \cdot \sim \alpha$, one may then conclude β , it can be shown that the theorems then derivable are precisely the theorems of the intuitionistic calculus.

THE PREDICATE CALCULUS

Propositions may also be built up, not out of other propositions but out of elements that are not themselves propositions. The simplest kind to be considered here are propositions in which a certain object or individual (in a wide sense) is said to possess a certain property or characteristic; e.g., "Socrates is wise" and "The number 7 is prime." Such a proposition contains two distinguishable parts: (1) an expression that names or designates an individual and (2) an expression, called a predicate, that stands for the property that that individual is said to possess. If x, y, z, \dots are used as individual variables (replaceable by names of individuals) and the symbols ϕ (phi), ψ (psi), χ (chi), \dots as predicate variables (replaceable by predicates), the formula ϕx is used to express the form of the propositions in question. Here x is said to be the argument of ϕ ; a predicate (or predicate variable) with only a single argument is said to be a monadic, or one-place, predicate (variable). Predicates with two or more arguments stand not for properties of single individuals but for relations between individuals. Thus the proposition "Tom is a son of John" is analyzable into two names of individuals ("Tom" and "John") and a dyadic or two-place predicate ("is a son of"), of which they are the arguments; and the proposition is thus of the form ϕxy . Analogously, "... is between ... and ..." is a three-place predicate, requiring three arguments, and so on. In general, a predicate variable followed by any number of individual variables is a wff of the predicate calculus. Such a wff is known as an atomic formula, and the predicate variable in it is said to be of degree n , if n is the number of individual variables following it. The degree of a predicate variable is sometimes indicated by a superscript—e.g., ϕ^3xyz may be

Fundamental definitions

written as ϕ^3xyz ; ϕ^3xy would then be regarded as not well formed. This practice is theoretically more accurate, but the superscripts are commonly omitted for ease of reading when no confusion is likely to arise.

Atomic formulas may be combined with truth-functional operators to give formulas such as $\phi x \vee \psi y$ [Example: "Either the customer (x) is friendly (ϕ) or else John (y) is disappointed (ψ)"]; $\phi xy \supset \sim \psi x$ [Example: "If the road (x) is above (ϕ) the flood line (y), then the road is not wet ($\sim \psi$)"]; and so on. Formulas so formed, however, are valid when and only when they are substitution-instances of valid wffs of PC and hence in a sense do not transcend PC. More interesting formulas are formed by the use, in addition, of quantifiers. There are two kinds of quantifiers: universal quantifiers, written as " $(\forall _)$ " or often simply as " $(_)$," where the blank is filled by a variable, which may be read "For all $_$ "; and existential quantifiers, written as " $(\exists _)$," which may be read "For some $_$ " or "There is a $_$ such that." ("Some" is to be understood as meaning "at least one.") Thus $(\forall x)\phi x$ is to mean "For all x , x is ϕ " or, more simply, "Everything is ϕ "; and $(\exists x)\phi x$ is to mean "For some x , x is ϕ " or, more simply, "Something is ϕ " or "There is a ϕ ." Slightly more complex examples are $(\forall x)(\phi x \supset \psi x)$ for "Whatever is ϕ is ψ ," $(\exists x)(\phi x \cdot \psi x)$ for "Something is both ϕ and ψ ," $(\forall x)(\exists y)\phi xy$ for "Everything bears the relation ϕ to at least one thing," and $(\exists x)(\forall y)\phi xy$ for "There is something that bears the relation ϕ to everything." To take a concrete case, if ϕxy means " x loves y " and the values of x and y are taken to be human beings, then the last two formulas mean, respectively, "Everybody loves somebody" and "Somebody loves everybody."

Quantifiers \forall and \exists

Intuitively, the notions expressed by the words "some" and "every" are connected in the following way: to assert that something has a certain property amounts to denying that everything lacks that property (for example, to say that something is white is to say that not everything is nonwhite); and, similarly, to assert that everything has a certain property amounts to denying that there is something that lacks it. These intuitive connections are reflected in the usual practice of taking one of the quantifiers as primitive and defining the other in terms of it. Thus \forall may be taken as primitive, and \exists introduced by the definition

$$(\exists a)\alpha =_{\text{Df}} \sim(\forall a)\sim\alpha$$

in which a is any variable and α is any wff; or, alternatively, \exists may be taken as primitive, and \forall introduced by the definition

$$(\forall a)\alpha =_{\text{Df}} \sim(\exists a)\sim\alpha.$$

The lower predicate calculus. A predicate calculus in which the only variables that occur in quantifiers are individual variables is known as a lower (or first-order) predicate calculus. Various lower predicate calculi have been constructed. In the most straightforward of these, to which the most attention will be devoted in this discussion and which subsequently will be referred to simply as LPC, the wffs can be specified as follows: Let the primitive symbols be (1) x, y, \dots (individual variables), (2) ϕ, ψ, \dots , each of some specified degree (predicate variables), and (3) the symbols $\sim, \vee, \forall, (,)$ and \cdot . An infinite number of each type of variable can now be secured as before by the use of numerical subscripts. The symbols \cdot, \supset , and \equiv are defined as in PC, and \exists as explained above. The formation rules are:

1. An expression consisting of a predicate variable of degree n followed by n individual variables is a wff.
2. If α is a wff, so is $\sim\alpha$.
3. If α and β are wffs, so is $(\alpha \vee \beta)$.
4. If α is a wff and a is an individual variable, then $(\forall a)\alpha$ is a wff. (In such a wff, α is said to be the scope of the quantifier.)

If a is any individual variable and α is any wff, every occurrence of a in α is said to be bound (by the quantifiers) when occurring in the wffs $(\forall a)\alpha$ and $(\exists a)\alpha$. Any occurrence of a variable that is not bound is said to be free. Thus, in $(\forall x)(\phi x \vee \psi y)$ the x in ϕx is bound, since it occurs within the scope of a quantifier containing x , but

Bound and free variables

y is free. In the wffs of a lower predicate calculus, every occurrence of a predicate variable (ϕ, ψ, χ, \dots) is free. A wff containing no free individual variables is said to be a closed wff of LPC. If a wff of LPC is considered as a proposition form, instances of it are obtained by replacing all free variables in it by predicates or by names of individuals, as appropriate. A bound variable, on the other hand, indicates not a point in the wff where a replacement is needed but a point (so to speak) at which the relevant quantifier applies.

For example, in ϕx , in which both variables are free, each variable must be replaced appropriately if a proposition of the form in question (such as "Socrates is white") is to be obtained; but in $(\exists x)\phi x$, in which x is bound, it is necessary only to replace ϕ by a predicate in order to obtain a complete proposition (e.g., replacing ϕ by "is white" yields the proposition "Something is white").

Validity in LPC. Intuitively, a wff of LPC is valid if and only if all its instances are true—i.e., if and only if every result of replacing each of its free variables appropriately and uniformly is a true proposition. A formal definition of validity in LPC to express this intuitive notion more precisely can be given as follows: for any wff of LPC, any number of LPC models can be formed. An LPC model has two elements: One is a set, D , of objects, known as a domain. D may contain as many or as few objects as one chooses, but it must contain at least one, and the objects may be of any kind. The other element, V , is a system of value assignments satisfying the following conditions. To each individual variable there is assigned some member of D (not necessarily a different one in each case). Assignments are next made to the predicate variables in the following way: if ϕ is monadic, there is assigned to it some subset of D (possibly the whole of D); intuitively this subset can be viewed as the set of all the objects in D that have the property ϕ . If ϕ is dyadic, there is assigned to it some set of ordered pairs (i.e., pairs of objects of which one is marked out as the first and the other as the second) drawn from D ; intuitively these can be viewed as all the pairs of objects in D in which the relation ϕ holds between the first object in the pair and the second. In general, if ϕ is of degree n , there is assigned to it some set of ordered n -tuples (groups of n objects) of members of D . It is then stipulated that an atomic formula is to have the value 1 in the model if the members of D assigned to its individual variables form, in that order, one of the n -tuples assigned to the predicate variable in it; otherwise, it is to have the value 0. Thus in the simplest case, ϕx will have the value 1 if the object assigned to x is one object in the set of objects assigned to ϕ ; and, if it is not, then ϕx will have the value 0. The values of truth functions are determined by the values of their arguments as in PC. Finally, the value of $(\forall x)\alpha$ is to be 1 if both (1) the value of α itself is 1 and (2) α would always still have the value 1 if a different assignment were made to x but all the other assignments were left precisely as they were; otherwise $(\forall x)\alpha$ is to have the value 0. Since \exists can be defined in terms of \forall , these rules cover all the wffs of LPC. A given wff may of course have the value 1 in some LPC models but the value 0 in others. But a valid wff of LPC may now be defined as one that has the value 1 in every LPC model. If 1 and 0 are viewed as representing truth and falsity, respectively, then validity is defined as truth in every model.

Although the above definition of validity in LPC is quite precise, it does not yield, as did the corresponding definition of PC validity in terms of truth tables, an effective decision procedure. It can, indeed, be shown that no generally applicable decision procedure for LPC is possible—i.e., that LPC is not a decidable system. This does not mean that it is never possible to prove that a given wff of LPC is valid—the validity of an unlimited number of such wffs can in fact be demonstrated—but it does mean that in the case of LPC, unlike that of PC, there is no general procedure, stated in advance, that would enable one to determine, for any wff whatever, whether it is valid or not (see also below *Model theory*).

Logical manipulations in LPC. The intuitive connections between *some* and *every* noted earlier are reflected in the fact that the following equivalences are valid:

$$(\exists x)\phi x \equiv \sim(\forall x)\sim\phi x$$

$$(\forall x)\phi x \equiv \sim(\exists x)\sim\phi x$$

These equivalences remain valid when ϕx is replaced by any wff, however complex; i.e., for any wff α whatsoever,

$$(\exists x)\alpha \equiv \sim(\forall x)\sim\alpha$$

and

$$(\forall x)\alpha \equiv \sim(\exists x)\sim\alpha$$

are valid. Because the rule of substitution of equivalents can be shown to hold in LPC, it follows that $(\exists x)$ may be replaced anywhere in a wff by $\sim(\forall x)\sim$, or $(\forall x)$ by $\sim(\exists x)\sim$, and the resulting wff will be equivalent to the original. Similarly, because the law of double negation permits the deletion of a pair of consecutive negation signs, $\sim(\exists x)$ may be replaced by $(\forall x)\sim$, and $\sim(\forall x)$ by $(\exists x)\sim$.

These principles are easily extended to more complex cases. To say that there is a pair of objects satisfying a certain condition is equivalent to denying that every pair of objects fails to satisfy that condition, and to say that every pair of objects satisfies a certain condition is equivalent to denying that there is any pair of objects that fails to satisfy that condition. These equivalences are expressed formally by the validity, again for any wff α , of

$$(\exists x)(\exists y)\alpha \equiv \sim(\forall x)(\forall y)\sim\alpha$$

and

$$(\forall x)(\forall y)\alpha \equiv \sim(\exists x)(\exists y)\sim\alpha$$

and by the resulting replaceability anywhere in a wff of $(\exists x)(\exists y)$ by $\sim(\forall x)(\forall y)\sim$, or of $(\forall x)(\forall y)$ by $\sim(\exists x)(\exists y)\sim$.

Analogously, $(\exists x)(\forall y)$ can be replaced by $\sim(\forall x)(\exists y)\sim$ [e.g., $(\exists x)(\forall y)(x \text{ loves } y)$ —"There is someone who loves everyone"—is equivalent to $\sim(\forall x)(\exists y)\sim(x \text{ loves } y)$ —"It is not true of everyone that there is someone whom he does not love"]; $(\forall x)(\exists y)$ can be replaced by $\sim(\exists x)(\forall y)\sim$; and in general the following rule, covering sequences of quantifiers of any length, holds:

1. If a wff contains an unbroken sequence of quantifiers, then the wff that results from replacing \forall by \exists and vice versa throughout that sequence and inserting or deleting \sim at each end of it is equivalent to the original wff.

This may be called the rule of quantifier transformation. It reflects, in a generalized form, the intuitive connections between "some" and "every" that were noted above.

The following are also valid, again where α is any wff:

$$(\forall x)(\forall y)\alpha \equiv (\forall y)(\forall x)\alpha$$

$$(\exists x)(\exists y)\alpha \equiv (\exists y)(\exists x)\alpha$$

The extensions of these lead to the following rule:

2. If a wff contains an unbroken sequence either of universal or of existential quantifiers, these quantifiers may be rearranged in any order and the resulting wff will be equivalent to the original wff.

This may be called the rule of quantifier rearrangement.

Two other important rules concern implications, not equivalences:

3. If a wff β begins with an unbroken sequence of quantifiers, and β' is obtained from β by replacing \forall by \exists at one or more places in the sequence, then β is stronger than β' —in the sense that $(\beta \supset \beta')$ is valid but $(\beta' \supset \beta)$ is in general not valid.
4. If a wff β begins with an unbroken sequence of quantifiers in which some existential quantifier Q_1 precedes some universal quantifier Q_2 , and if β' is obtained from β by moving Q_1 to the right of Q_2 , then β is stronger than β' .

As illustrations of these rules, the following are valid for any wff α :

$(\forall x)(\forall y)\alpha \supset (\exists x)(\forall y)\alpha$	rule 3
$(\exists x)(\forall y)(\forall z)\alpha \supset (\exists x)(\forall y)(\exists z)\alpha$	rule 3
$(\exists x)(\forall y)\alpha \supset (\forall y)(\exists x)\alpha$	rule 4
$(\exists x)(\exists y)(\forall z)\alpha \supset (\exists y)(\forall z)(\exists x)\alpha$	rule 4

In each case the converses are not valid (though they may be valid in particular cases in which α is of some special form).

Value assignments in the domain of a model

Laws concerning quantifiers

Application of the rules

Some of the uses of the above rules can be illustrated by considering a wff α that contains precisely two free individual variables. By prefixing to α two appropriate quantifiers and possibly one or more negation signs, it is possible to form a closed wff (called a closure of α) that will express a determinate proposition when a meaning is assigned to the predicate variables. The above rules can be used to list exhaustively the nonequivalent closures of α and the implication relations between them. The simplest example is ϕxy , which for illustrative purposes can be taken to mean "x loves y." Application of rules 1 and 2 will show that every closure of ϕxy is equivalent to one or another of the following 12 wffs (none of which is in fact equivalent to any of the others):

- (a) $(\forall x)(\forall y)\phi xy$ ("Everybody loves everybody");
- (b) $(\exists x)(\forall y)\phi xy$ ("Somebody loves everybody");
- (c) $(\exists y)(\forall x)\phi xy$ ("There is someone whom everyone loves");
- (d) $(\forall y)(\exists x)\phi xy$ ("Each person is loved by at least one person");
- (e) $(\forall x)(\exists y)\phi xy$ ("Each person loves at least one person");
- (f) $(\exists x)(\exists y)\phi xy$ ("Somebody loves somebody"); and
- (g)-(l) the respective negations of each of the above.

Rules 3 and 4 show that the following implications among formulas (a)-(f) are valid:

$$\begin{array}{lll} (a) \supset (b) & (d) \supset (f) & (c) \supset (e) \\ (b) \supset (d) & (a) \supset (c) & (e) \supset (f) \end{array}$$

The implications holding among the negations of (a)-(f) follow from these by the law of transposition; e.g., since $(a) \supset (b)$ is valid, so is $\sim(b) \supset \sim(a)$. The quantification of wffs containing three, four, etc., variables can be dealt with by the same rules.

Intuitively, $(\forall x)\phi x$ and $(\forall y)\phi y$ both "say the same thing"—namely, that everything is ϕ —and $(\exists x)\phi x$ and $(\exists y)\phi y$ both mean simply that something is ϕ . Clearly, so long as the same variable occurs both in the quantifier and as the argument of ϕ , it does not matter what letter is chosen for this purpose. The procedure of replacing some variable in a quantifier, together with every occurrence of that variable in its scope, by some other variable that does not occur elsewhere in its scope is known as relettering a bound variable. If β is the result of relettering a bound variable in a wff α , then α and β are said to be bound alphabetical variants of each other, and bound alphabetical variants are always equivalent. The reason for restricting the replacement variable to one not occurring elsewhere in the scope of the quantifier can be seen from an example: If ϕxy is taken as before to mean "x loves y," the wff $(\forall x)\phi xy$ expresses the proposition form "Everyone loves y," in which the identity of y is left unspecified, and so does its bound alphabetical variant $(\forall z)\phi zy$. If x were replaced by y , however, the closed wff $(\forall y)\phi yy$ would be obtained, which expresses the proposition that everyone loves himself and is clearly not equivalent to the original.

A wff in which all the quantifiers occur in an unbroken sequence at the beginning, with the scope of each extending to the end of the wff, is said to be in prenex normal form (PNF). Wffs that are in PNF are often more convenient to work with than those that are not. For every wff of LPC, however, there is an equivalent wff in PNF (often simply called its PNF). One effective method for finding the PNF of any given wff is the following:

1. Reletter bound variables as far as is necessary to ensure (a) that each quantifier contains a distinct variable and (b) that no variable in the wff occurs both bound and free.
2. Use definitions or PC equivalences to eliminate all operators except \sim , \cdot , and \vee .
3. Use the De Morgan laws and the rule of quantifier transformation to eliminate all occurrences of \sim immediately before parentheses or quantifiers.
4. Gather all of the quantifiers into a sequence at the beginning in the order in which they appear in the wff and take the whole of what remains as their scope. Example:

$$(\forall x)[(\phi x \cdot (\exists y)\psi xy) \supset (\exists y)\chi xy] \supset (\exists z)(\phi z \supset \psi zx).$$

Step 1 can be achieved by relettering the third and fourth occurrences of y and every occurrence of x except the last (which is free); thus

$$(\forall w)[(\phi w \cdot (\exists y)\psi wy) \supset (\exists u)\chi wu] \supset (\exists z)(\phi z \supset \psi zx).$$

Step 2 now yields

$$\sim(\forall w)(\sim[\phi w \cdot (\exists y)\psi wy] \vee (\exists u)\chi wu) \vee (\exists z)(\sim\phi z \vee \psi zx).$$

By step 3 this becomes

$$(\exists w)[(\phi w \cdot (\exists y)\psi wy) \cdot (\forall u)\sim\chi wu] \vee (\exists z)(\sim\phi z \vee \psi zx).$$

Finally, step 4 yields

$$(\exists w)(\exists y)(\forall u)(\exists z)[(\phi w \cdot \psi wy) \cdot \sim\chi wu] \vee (\sim\phi z \vee \psi zx),$$

which is in PNF.

Classification of dyadic relations. Consider the closed wff,

$$(\forall x)(\forall y)(\phi xy \supset \phi yx),$$

which means that, whenever the relation ϕ holds between one object and a second, it also holds between that second object and the first. This expression is not valid, since it is true for some relations but false for others. A relation for which it is true is called a symmetrical relation (example: "is parallel to"). If the relation ϕ is such that, whenever it holds between one object and a second, it fails to hold between the second and the first—i.e., if ϕ is such that

$$(\forall x)(\forall y)(\phi xy \supset \sim\phi yx)$$

—then ϕ is said to be asymmetrical (example: "is greater than"). A relation that is neither symmetrical nor asymmetrical is said to be nonsymmetrical. Thus ϕ is nonsymmetrical if

$$(\exists x)(\exists y)(\phi xy \cdot \phi yx) \cdot (\exists x)(\exists y)(\phi xy \cdot \sim\phi yx)$$

(example: "loves").

Dyadic relations can also be characterized in terms of another threefold division: A relation ϕ is said to be transitive if, whenever it holds between one object and a second and also between that second object and a third, it holds between the first and the third—i.e., if

$$(\forall x)(\forall y)(\forall z)[(\phi xy \cdot \phi yz) \supset \phi xz]$$

(example: "is greater than"). An intransitive relation is one that, whenever it holds between one object and a second and also between that second and a third, fails to hold between the first and the third; i.e., ϕ is intransitive if

$$(\forall x)(\forall y)(\forall z)[(\phi xy \cdot \phi yz) \supset \sim\phi xz]$$

(example: "is father of"). A relation that is neither transitive nor intransitive is said to be nontransitive. Thus ϕ is nontransitive if

$$(\exists x)(\exists y)(\exists z)(\phi xy \cdot \phi yz \cdot \phi xz) \cdot (\exists x)(\exists y)(\exists z)(\phi xy \cdot \phi yz \cdot \sim\phi xz)$$

(example: "is a first cousin of").

A relation ϕ that always holds between any object and itself is said to be reflexive; i.e., ϕ is reflexive if

$$(\forall x)\phi xx$$

(example: "is identical with"). If ϕ never holds between any object and itself—i.e., if

$$\sim(\exists x)\phi xx$$

—then ϕ is said to be irreflexive (example: "is greater than"). If ϕ is neither reflexive nor irreflexive—i.e., if

$$(\exists x)\phi xx \cdot (\exists x)\sim\phi xx$$

—then ϕ is said to be nonreflexive (example: "admires").

A relation such as "is of the same length as" is not strictly reflexive, as some objects do not have a length at all and thus are not of the same length as anything, even themselves. But this relation is reflexive in the weaker sense that whenever an object is of the same length as anything it is of the same length as itself. Such a relation is said to be quasi-reflexive. Thus ϕ is quasi-reflexive if

$$(\forall x)[(\exists y)\phi xy \supset \phi xx].$$

A reflexive relation is of course also quasi-reflexive.

Symmetry, transitivity, reflexivity

Prenex normal form

For the most part, these three classifications are independent of each other; thus a symmetrical relation may be transitive (like "is equal to") or intransitive (like "is perpendicular to") or nontransitive (like "is one mile distant from"). There are, however, certain limiting principles, of which the most important are:

1. Every relation that is symmetrical and transitive is at least quasi-reflexive.
2. Every asymmetrical relation is irreflexive.
3. Every relation that is transitive and irreflexive is asymmetrical.

A relation that is reflexive, symmetrical, and transitive is called an equivalence relation.

Axiomatization of LPC. Rules of uniform substitution for predicate calculi, though formulable, are mostly very complicated, and, to avoid the necessity for these rules, axioms for these systems are therefore usually given by axiom schemata in the sense explained earlier in *Axiomatization of PC*. Given the formation rules and definitions stated in the introductory paragraph of *The lower predicate calculus*, the following is presented as one standard axiomatic basis for LPC:

Axiom schemata:

1. Any LPC substitution-instance of any valid wff of PC is an axiom.
2. Any wff of the form $(\forall a)\alpha \supset \beta$ is an axiom, if β is either identical with α or differs from it only in that, wherever α has a free occurrence of a , β has a free occurrence of some other individual variable b .
3. Any wff of the form $(\forall a)(\alpha \supset \beta) \supset [\alpha \supset (\forall a)\beta]$ is an axiom, provided that α contains no free occurrence of a .

Transformation rules:

1. Modus ponens (as given above in *Axiomatization of PC*).
2. If α is a theorem, so is $(\forall a)\alpha$, where a is any individual variable (rule of universal generalization).

The axiom schemata call for some explanation and comment. By an LPC substitution-instance of a wff of PC is meant any result of uniformly replacing every propositional variable in that wff by a wff of LPC. Thus one LPC substitution-instance of $(p \supset \sim q) \supset (q \supset \sim p)$ is $[\phi xy \supset \sim(\forall x)\psi x] \supset [(\forall x)\psi x \supset \sim\phi xy]$. Axiom schema 1 makes available in LPC all manipulations such as commutation, transposition, and distribution, which depend only on PC principles. Examples of wffs that are axioms by axiom schema 2 are $(\forall x)\phi x \supset \phi x$, $(\forall x)\phi x \supset \phi y$, and $(\forall x)(\exists y)\phi xy \supset (\exists y)\phi zy$. To see why it is necessary for the variable that replaces a to be free in β , consider the last example: Here a is x , α is $(\exists y)\phi xy$, in which x is free, and β is $(\exists y)\phi zy$, in which z is free and replaces x . But had y , which would become bound by the quantifier $(\exists y)$, been chosen as a replacement instead of z , the result would have been $(\forall x)(\exists y)\phi xy \supset (\exists y)\phi yy$, the invalidity of which can be seen intuitively by taking ϕxy to mean "x is a child of y": for then $(\forall x)(\exists y)\phi xy$ will mean that everyone is a child of someone, which is true, but $(\exists y)\phi yy$ will mean that someone is a child of himself, which is false. The need for the proviso in axiom schema 3 can also be seen from an example. Defiance of the proviso would give as an axiom $(\forall x)(\phi x \supset \psi x) \supset [\phi x \supset (\forall x)\psi x]$; and, if ϕx were taken to mean "x is a Spaniard," ψx to mean "x is a European," and the free occurrence of x (the first occurrence in the consequent) to stand for General Franco, then the antecedent would mean that every Spaniard is a European, but the consequent would mean that, if General Franco is a Spaniard, then everyone is a European.

It can be proved—though the proof is not an elementary one—that the theorems derivable from the above basis are precisely the wffs of LPC that are valid by the definition of validity given above in *Validity in LPC*. Several other bases for LPC are known that also have this property. The axiom schemata and transformation rules here given are such that any purported proof of a theorem can be effectively checked to determine whether it really is a proof or not; nevertheless, theoremhood in LPC, like validity in LPC, is not effectively decidable, in that there is no effective method of telling with regard to any arbitrary wff

whether it is a theorem or not. In this respect, axiomatic bases for LPC contrast with those for PC.

Semantic tableaux. Since the 1980s another technique for determining the validity of arguments in either PC or LPC has gained some popularity, owing both to its ease of learning and to its straightforward implementation by computer programs. Originally suggested by the Dutch logician Evert W. Beth (1908–64), it was more fully developed and publicized by the American mathematician and logician Raymond M. Smullyan (b. 1919). Resting on the observation that it is impossible for the premises of a valid argument to be true while the conclusion is false, this method attempts to interpret (or evaluate) the premises in such a way that they are all simultaneously satisfied and the negation of the conclusion is also satisfied. Success in such an effort would show the argument to be invalid, while failure to find such an interpretation would show it to be valid.

The construction of a semantic tableau proceeds as follows: Express the premises and negation of the conclusion of an argument in PC using only negation (\sim) and disjunction (\vee) as propositional connectives. Eliminate every occurrence of two negation signs in a sequence (e.g., $\sim\sim\sim\sim a$ becomes $\sim a$). Now construct a tree diagram branching downward such that each disjunction is replaced by two branches, one for the left disjunct and one for the right. The original disjunction is true if either branch is true. Reference to De Morgan's laws shows that a negation of a disjunction is true just in case the negations of both disjuncts are true (i.e., $\sim(p \vee q) \equiv (\sim p \cdot \sim q)$). This semantic observation leads to the rule that the negation of a disjunction becomes one branch containing the negation of each disjunct:

$$\begin{array}{c} \sim(a \vee b) \\ | \\ \sim a \\ \sim b \end{array}$$

Consider the following argument:

$$\frac{a \vee b}{\frac{\sim a}{b}}$$

Write:

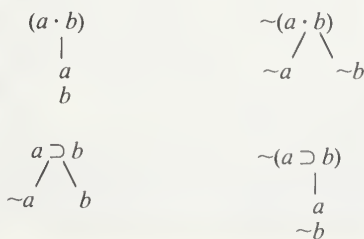
$$\begin{array}{c} a \vee b \\ \sim a \\ \sim b \end{array}$$

Now strike out the disjunction and form two branches:



Only if all the sentences in at least one branch are true is it possible for the original premises to be true and the conclusion false (equivalently for the negation of the conclusion). By tracing the line upward in each branch to the top of the tree, one observes that no valuation of a in the left branch will result in all the sentences in that branch receiving the value true (because of the presence of a and $\sim a$). Similarly, in the right branch the presence of b and $\sim b$ makes it impossible for a valuation to result in all the sentences of the branch receiving the value true. These are all the possible branches; thus, it is impossible to find a situation in which the premises are true and the conclusion false. The original argument is therefore valid.

This technique can be extended to deal with other connectives:



Axiomatic basis and undecidability of theorems

Constructing a tableau

Furthermore, in LPC, rules for instantiating quantified wffs need to be introduced. Clearly, any branch containing both $(\forall x)\phi x$ and $\sim\phi y$ is one in which not all the sentences in that branch can be simultaneously satisfied (under the assumption of ω -consistency; see below *Metalogic*). Again, if all the branches fail to be simultaneously satisfiable, the original argument is valid.

Special systems of LPC. LPC as expounded above may be modified by either restricting or extending the range of wffs in various ways.

1. Partial systems of LPC. Some of the more important systems produced by restriction are here outlined:

a. It may be required that every predicate variable be monadic, while still allowing an infinite number of individual and predicate variables. The atomic wffs are then simply those consisting of a predicate variable followed by a single individual variable. Otherwise, the formation rules remain as before, and the definition of validity is also as before, though simplified in obvious ways. This system is known as the monadic LPC; it provides a logic of properties but not of relations. One important characteristic of this system is that it is decidable. (The introduction of even a single dyadic predicate variable, however, would make the system undecidable, and, in fact, even the system that contains only a single dyadic predicate variable and no other predicate variables at all has been shown to be undecidable.)

b. A still simpler system can be formed by requiring (1) that every predicate variable be monadic, (2) that only a single individual variable (e.g., x) be used, (3) that every occurrence of this variable be bound, and (4) that no quantifier occur within the scope of any other. Examples of wffs of this system are $(\forall x)(\phi x \supset (\psi x \cdot \chi x))$ ("Whatever is ϕ is both ψ and χ "); $(\exists x)(\phi x \cdot \sim\psi x)$ ("There is something that is ϕ but not ψ "); and $(\forall x)(\phi x \supset \psi x) \supset (\exists x)(\phi x \cdot \psi x)$ ("If whatever is ϕ is ψ , then something is both ϕ and ψ "). The notation for this system can be simplified by omitting x everywhere and writing $\exists\phi$ for "Something is ϕ ," $\forall(\phi \supset \psi)$ for "Whatever is ϕ is ψ ," and so on. Although this system is more rudimentary even than the monadic LPC (of which it is a fragment), the forms of a wide range of inferences can be represented in it. It is also a decidable system, and decision procedures of an elementary kind can be given for it.

2. Extensions of LPC. More elaborate systems, in which a wider range of propositions can be expressed, have been constructed, by adding to LPC new symbols of various types. The most straightforward of such additions are:

a. One or more individual constants (say a, b, \dots): these constants are interpreted as names of specific individuals; formally they are distinguished from individual variables by the fact that they cannot occur within quantifiers: e.g., $(\forall x)$ is a quantifier but $(\forall a)$ is not.

b. One or more predicate constants (say A, B, \dots), each of some specified degree, thought of as designating specific properties or relations.

A further possible addition, which calls for somewhat fuller explanation, consists of symbols designed to stand for functions. The notion of a function may be sufficiently explained for present purposes as follows: There is said to be a certain function of n arguments (or, of degree n) when there is a rule that specifies a unique object (called the value of the function) whenever all the arguments are specified. In the domain of human beings, for example, "the mother of—" is a monadic function (a function of one argument), since for every human being there is a unique individual who is his mother; and in the domain of the natural numbers (i.e., 0, 1, 2, . . .), "the sum of— and —" is a function of two arguments, since for any pair of natural numbers there is a natural number that is their sum. A function symbol can be thought of as forming a name out of other names (its arguments): thus, whenever x and y name numbers, "the sum of x and y " also names

a number, and similarly for other kinds of functions and arguments.

To enable functions to be expressed in LPC there may be added:

c. One or more function variables (say f, g, \dots) or one or more function constants (say, F, G, \dots) or both, each of some specified degree. The former are interpreted as ranging over functions of the degrees specified, and the latter as designating specific functions of that degree.

When any or all of $a-c$ are added to LPC, the formation rules listed in the first paragraph of the section *The lower predicate calculus* need to be modified to enable the new symbols to be incorporated into wffs. This can be done as follows: A term is first defined as either (1) an individual variable or (2) an individual constant or (3) any expression formed by prefixing a function variable or function constant of degree n to any n terms (these terms—the arguments of the function symbol—are usually separated by commas and enclosed in parentheses). Formation rule 1 is then replaced by:

1'. An expression consisting of a predicate variable or predicate constant of degree n followed by n terms is a wff.

The axiomatic basis given in *Axiomatization of LPC* also requires the following modification: in axiom schema 2 any term is allowed to replace a when β is formed, provided that no variable that is free in the term becomes bound in β . The following examples will illustrate the use of the aforementioned additions to LPC: Let the values of the individual variables be the natural numbers; let the individual constants a and b stand for the numbers 2 and 3, respectively; let A mean "is prime"; and let F represent the dyadic function "the sum of." Then $AF(a,b)$ expresses the proposition "The sum of 2 and 3 is prime," and $(\exists x) AF(x,a)$ expresses the proposition "There exists a number such that the sum of it and 2 is prime."

The introduction of constants is normally accompanied by the addition, to the axiomatic basis, of special axioms containing those constants, designed to express principles that hold of the objects, properties, relations, or functions represented by them—though they do not hold of objects, properties, relations, or functions in general. It may be decided, for example, to use the constant A to represent the dyadic relation "is greater than" (so that Axy is to mean " x is greater than y ," and so forth). This relation, unlike many others, is transitive; i.e., if one object is greater than a second and that second is in turn greater than a third, then the first is greater than the third. Hence, the following special axiom schema might be added: If $t_1, t_2,$ and t_3 are any terms, then

$$(At_1t_2 \cdot At_2t_3) \supset At_1t_3$$

is an axiom. By such means systems can be constructed to express the logical structures of various particular disciplines. The area in which most work of this kind has been done is that of natural-number arithmetic (see below *Axiomatization of arithmetic*).

PC and LPC are sometimes combined into a single system. This may be done most simply by adding propositional variables to the list of LPC primitives, adding a formation rule to the effect that a propositional variable standing alone is a wff, and deleting "LPC" in axiom schema 1. This yields as wffs such expressions as $(p \vee q) \supset (\forall x)\phi x$ and $(\exists x)[p \supset (\forall y)\phi xy]$.

3. LPC-with-identity. The word "is" is not always used in the same way. In a proposition such as (1) "Socrates is snub-nosed," the expression preceding the "is" names an individual and the expression following it stands for a property attributed to that individual. But, in a proposition such as (2) "Socrates is the Athenian philosopher who drank hemlock," the expressions preceding and following the "is" both name individuals, and the sense of the whole proposition is that the individual named by the first is the same individual as the individual named by the second. Thus in (2) "is" can be expanded to "is the same individual as," whereas in (1) it cannot. As used in (2), "is" stands for a dyadic relation—namely, identity—that the propo-

Monadic
LPC and
subsystem

Special
axioms

Defini-
tions:
"function,"
"degree,"
"term"

Identity propositions

sition asserts to hold between the two individuals. An identity proposition is to be understood in this context as asserting no more than this; in particular it is not to be taken as asserting that the two naming expressions have the same meaning. A much discussed example to illustrate this last point is "The morning star is the evening star." It is false that the expressions "the morning star" and "the evening star" mean the same, but true that the object referred to by the former is the same as that referred to by the latter.

To enable the forms of identity propositions to be expressed, a dyadic predicate constant is added to LPC, for which the most usual notation is = (written between, rather than before, its arguments). The intended interpretation of $x = y$ is that x is the same individual as y , and the most convenient reading is "x is identical with y." Its negation $\neg(x = y)$ is commonly abbreviated to $x \neq y$. To the definition of an LPC model in *Validity in LPC* there is now added the rule (which accords in an obvious way with the intended interpretation) that the value of $x = y$ is to be 1 if the same member of D is assigned to both x and y and that otherwise its value is to be 0; validity can then be defined as before. The following additions (or some equivalent ones) are made to the axiomatic basis for LPC: the axiom (I1) $x = x$ and the axiom schema (I2) that, where α and β are any individual variables and α and β are wffs that differ only in that, at one or more places where α has a free occurrence of a , β has a free occurrence of b , $(\alpha = b) \supset (\alpha \supset \beta)$ is an axiom. Such a system is known as a lower-predicate-calculus-with-identity; it may of course be further augmented in the other ways referred to above in *Extensions of LPC*, in which case any term may be an argument of =.

Identity is an equivalence relation; i.e., it is reflexive, symmetrical, and transitive. Its reflexivity is directly expressed in axiom I1, and theorems expressing its symmetry and transitivity can easily be derived from the basis given.

Certain wffs of LPC-with-identity express propositions about the number of things that possess a given property. "At least one thing is ϕ " could, of course, already be expressed by $(\exists x)\phi x$; "At least two distinct (nonidentical) things are ϕ " can now be expressed by $(\exists x)(\exists y)(\phi x \cdot \phi y \cdot x \neq y)$; and the sequence can be continued in an obvious way. "At most one thing is ϕ " (i.e., "No two distinct things are both ϕ ") can be expressed by the negation of the last mentioned wff or by its equivalent, $(\forall x)(\forall y)[(\phi x \cdot \phi y) \supset x = y]$, and the sequence can again be easily continued. A formula for "Exactly one thing is ϕ " may be obtained by conjoining the formulas for "At least one thing is ϕ " and "At most one thing is ϕ ," but a simpler wff equivalent to this conjunction is $(\exists x)[\phi x \cdot (\forall y)(\phi y \supset x = y)]$, which means "There is something that is ϕ , and anything that is ϕ is that thing." The proposition "Exactly two things are ϕ " can be represented by

$$(\exists x)(\exists y)(\phi x \cdot \phi y \cdot x \neq y \cdot (\forall z)[\phi z \supset (z = x \vee z = y)]);$$

i.e., "There are two nonidentical things each of which is ϕ , and anything that is ϕ is one or the other of these." Clearly, this sequence can also be extended to give a formula for "Exactly n things are ϕ " for every natural number n . It is convenient to abbreviate the wff for "Exactly one thing is ϕ " to $(\exists!x)\phi x$. This special quantifier is frequently read aloud as "E-Shriek x ."

When a certain property ϕ belongs to one and only one object, it is convenient to have an expression that names that object. A common notation for this purpose is $(\iota x)\phi x$, which may be read as "the thing that is ϕ " or more briefly as "the ϕ ." In general, where a is any individual variable and α is any wff, $(\iota a)\alpha$ then stands for the single value of a that makes α true. An expression of the form "the so-and-so" is called a definite description; and (ιx) , known as a description operator, can be thought of as forming a name of an individual out of a proposition form. (ιx) is analogous to a quantifier in that, when prefixed to a wff α , it binds every free occurrence of x in α . Relettering of bound variables is also permissible: in the simplest case, $(\iota x)\phi x$ and $(\iota y)\phi y$ can each be read simply as "the ϕ ."

As far as formation rules are concerned, definite descrip-

tions can be incorporated into LPC by letting expressions of the form $(\iota a)\alpha$ count as terms; rule 1' of *Extensions of LPC* will then allow them to occur in atomic formulas (including identity formulas). "The ϕ is (i.e., has the property) ψ " can then be expressed as $\psi(\iota x)\phi x$; "y is (the same individual as) the ϕ " as $y = (\iota x)\phi x$; "The ϕ is (the same individual as) the ψ " as $(\iota x)\phi x = (\iota y)\psi y$; and so forth. The correct analysis of propositions containing definite descriptions has been the subject of considerable philosophical controversy. One widely accepted account, however—substantially that presented in *Principia Mathematica* and known as Russell's theory of descriptions, after Bertrand Russell—holds that "The ϕ is ψ " is to be understood as meaning that exactly one thing is ϕ and that thing is also ψ . In that case it can be expressed by a wff of LPC-with-identity that contains no description operators—namely,

$$(\exists x)[\phi x \cdot (\forall y)(\phi y \supset x = y) \cdot \psi x]. \tag{1}$$

Analogously, "y is the ϕ " is analyzed as "y is ϕ and nothing else is ϕ ," and hence as expressible by

$$\phi y \cdot (\forall x)(\phi x \supset x = y); \tag{2}$$

and "The ϕ is the ψ " is analyzed as "Exactly one thing is ϕ , exactly one thing is ψ , and whatever is ϕ is ψ ," and hence as expressible by

$$(\exists x)[\phi x \cdot (\forall y)(\phi y \supset x = y)] \cdot (\exists x)[\psi x \cdot (\forall y)(\psi y \supset x = y)] \cdot (\forall x)(\phi x \supset \psi x) \tag{3}$$

$\psi(\iota x)\phi x$, $y = (\iota x)\phi x$ and $(\iota x)\phi x = (\iota y)\psi y$ can then be regarded as abbreviations for (1), (2), and (3), respectively; and by generalizing to more complex cases, all wffs that contain description operators can be regarded as abbreviations for longer wffs that do not.

The analysis that leads to (1) as a formula for "The ϕ is ψ " leads to the following for "The ϕ is not ψ ":

$$(\exists x)[\phi x \cdot (\forall y)(\phi y \supset x = y) \cdot \sim \psi x]. \tag{4}$$

It is important to note that (4) is not the negation of (1); this negation is, instead,

$$\sim(\exists x)[\phi x \cdot (\forall y)(\phi y \supset x = y) \cdot \psi x]. \tag{5}$$

The difference in meaning between (4) and (5) lies in the fact that (4) is true only when there is exactly one thing that is ϕ and that thing is not ψ , but (5) is true both in this case and also when nothing is ϕ at all and when more than one thing is ϕ . Neglect of the distinction between (4) and (5) can result in serious confusion of thought: in ordinary speech it is frequently unclear whether someone who denies that the ϕ is ψ is conceding that exactly one thing is ϕ but denying that it is ψ , or denying that exactly one thing is ϕ .

The basic contention of Russell's theory of descriptions is that a proposition containing a definite description is not to be regarded as an assertion about an object of which that description is a name but rather as an existentially quantified assertion that a certain (rather complex) property has an instance. Formally, this is reflected in the rules for eliminating description operators that were outlined above.

Higher-order predicate calculi. A feature shared by LPC and all its extensions so far mentioned is that the only variables that occur in quantifiers are individual variables. It is in virtue of this feature that they are called lower (or first-order) calculi. Various predicate calculi of higher order can be formed, however, in which quantifiers may contain other variables as well, hence binding all free occurrences of these that lie within their scope. In particular, in the second-order predicate calculus, quantification is permitted over both individual and predicate variables; hence wffs such as $(\forall \phi)(\exists x)\phi x$ can be formed. This last formula, since it contains no free variables of any kind, expresses a determinate proposition—namely, the proposition that every property has at least one instance. One important feature of this system is that in it identity need not be taken as primitive but can be introduced by defining $x = y$ as $(\forall \phi)(\phi x \equiv \phi y)$ —i.e., "Every property possessed by x is also possessed by y and vice versa." Whether such a definition is acceptable as a general account of identity

Negations with differing scopes

Numerical quantification

Binding of predicate variables

Definite descriptions

is a question that raises philosophical issues too complex to be discussed here; they are substantially those raised by the principle of the identity of indiscernibles, best known for its exposition in the 17th century by Leibniz.

(G.E.H./M.L.Sc.)

MODAL LOGIC

True propositions can be divided into those—like “ $2 + 2 = 4$ ”—that are true by logical necessity (necessary propositions), and those—like “France is a republic”—that are not (contingently true propositions). Similarly, false propositions can be divided into those—like “ $2 + 2 = 5$ ”—that are false by logical necessity (impossible propositions), and those—like “France is a monarchy”—that are not (contingently false propositions). Contingently true and contingently false propositions are known collectively as contingent propositions. A proposition that is not impossible (*i.e.*, one that is either necessary or contingent) is said to be a possible proposition. Intuitively, the notions of necessity and possibility are connected in the following way: to say that a proposition is necessary is to say that it is not possible for it to be false; and to say that a proposition is possible is to say that it is not necessarily false.

If it is logically impossible for a certain proposition, p , to be true without a certain proposition, q , being also true (*i.e.*, if the conjunction of p and not- q is logically impossible), then it is said that p strictly implies q . An alternative, equivalent way of explaining the notion of strict implication is by saying that p strictly implies q if and only if it is necessary that p materially implies q . “John’s tie is scarlet,” for example, strictly implies “John’s tie is red,” because it is impossible for John’s tie to be scarlet without being red (or: it is necessarily true that if John’s tie is scarlet it is red); and in general, if p is the conjunction of the premises, and q the conclusion, of a deductively valid inference, p will strictly imply q .

The notions just referred to—necessity, possibility, impossibility, contingency, strict implication—and certain other closely related ones are known as modal notions, and a logic designed to express principles involving them is called a modal logic.

The most straightforward way of constructing such a logic is to add to some standard nonmodal system a new primitive operator intended to represent one of the modal notions mentioned above, to define other modal operators in terms of it, and to add certain special axioms or transformation rules or both. A great many systems of modal logic have been constructed, but attention will be restricted here to a few closely related ones in which the underlying nonmodal system is ordinary PC.

Alternative systems of modal logic. All the systems to be considered here have the same wffs but differ in their axioms. The wffs can be specified by adding to the symbols of PC a primitive monadic operator L and to the formation rules of PC the rule that if α is a wff, so is $L\alpha$. L is intended to be interpreted as “It is necessary that,” so that Lp will be true if and only if p is a necessary proposition. The monadic operator M and the dyadic operator \rightarrow (to be interpreted as “It is possible that” and “strictly implies,” respectively) can then be introduced by the following definitions, which reflect in an obvious way the informal accounts given above of the connections between necessity, possibility, and strict implication: if α is any wff, then $M\alpha$ is to be an abbreviation of $\sim L\sim\alpha$; and if α and β are any wffs, then $\alpha \rightarrow \beta$ is to be an abbreviation of $L(\alpha \supset \beta)$ [or alternatively of $\sim M(\alpha \cdot \sim\beta)$].

The modal system known as **T** has as axioms some set of axioms adequate for PC (such as those of PM), and in addition

1. $Lp \supset p$
2. $L(p \supset q) \supset (Lp \supset Lq)$

Axiom 1 expresses the principle that whatever is necessarily true is true, and 2 the principle that, if q logically follows from p , then, if p is a necessary truth, so is q (*i.e.*, that whatever follows from a necessary truth is itself a necessary truth). These two principles seem to have a high degree of intuitive plausibility, and 1 and 2 are theorems in almost all modal systems. The transformation rules of **T** are uniform substitution, modus ponens, and a rule to

the effect that if α is a theorem so is $L\alpha$ (the rule of necessitation). The intuitive rationale of this rule is that, in a sound axiomatic system, it is expected that every instance of a theorem α will be not merely true but necessarily true—and in that case every instance of $L\alpha$ will be true. Among the simpler theorems of **T** are

$$\begin{aligned} p &\supset Mp, \\ L(p \cdot q) &\equiv (Lp \cdot Lq), \\ M(p \vee q) &\equiv (Mp \vee Mq), \\ (Lp \vee Lq) &\supset L(p \vee q) \text{ (but not its converse),} \\ M(p \cdot q) &\supset (Mp \cdot Mq) \text{ (but not its converse),} \end{aligned}$$

and

$$\begin{aligned} LMp &\equiv \sim ML\sim p, \\ (p \rightarrow q) &\supset (Mp \supset Mq), \\ (\sim p \rightarrow p) &\equiv Lp, \\ L(p \vee q) &\supset (Lp \vee Mq). \end{aligned}$$

There are many modal formulas that are not theorems of **T** but that have a certain claim to express truths about necessity and possibility. Among them are

$$Lp \supset LLp, Mp \supset LMp, \text{ and } p \supset LMp.$$

The first of these means that if a proposition is necessary, its being necessary is itself a necessary truth; the second means that if a proposition is possible, its being possible is a necessary truth; and the third means that if a proposition is true, then not merely is it possible but its being possible is a necessary truth. These are all various elements in the general thesis that a proposition’s having the modal characteristics it has (such as necessity, possibility) is not a contingent matter but is determined by logical considerations. Although this thesis may be philosophically controversial, it is at least plausible, and its consequences are worth exploring. One way of exploring them is to construct modal systems in which the formulas listed above are theorems. None of these formulas, as was said, is a theorem of **T**; but each could be consistently added to **T** as an extra axiom to produce a new and more extensive system. The system obtained by adding $Lp \supset LLp$ to **T** is known as **S4**; that obtained by adding $Mp \supset LMp$ to **T** is known as **S5**; and the addition of $p \supset LMp$ to **T** gives the Brouwerian system, here called **B** for short.

The systems **S4**, **S5**, and **B**

The relations between these four systems are as follows: **S4** is stronger than **T**—*i.e.*, it contains all the theorems of **T** and others besides. **B** is also stronger than **T**. **S5** is stronger than **S4** and also stronger than **B**. **S4** and **B**, however, are independent of each other in the sense that each contains some theorems that the other does not have. It is of particular importance that if $Mp \supset LMp$ is added to **T** then $Lp \supset LLp$ can be derived as a theorem but that if one merely adds the latter to **T** the former cannot then be derived.

Examples of theorems of **S4** that are not theorems of **T** are $Mp \equiv MMp$, $MLMp \supset Mp$, and $(p \rightarrow q) \supset (Lp \rightarrow Lq)$. Examples of theorems of **S5** that are not theorems of **S4** are $Lp \equiv MLp$, $L(p \vee Mq) \equiv (Lp \vee Mq)$, $M(p \cdot Lq) \equiv (Mp \cdot Lq)$, and $(Lp \rightarrow Lq) \vee (Lq \rightarrow Lp)$. One important feature of **S5** but not of the other systems mentioned is that any wff that contains an unbroken sequence of monadic modal operators (L s or M s or both) is provably equivalent to the same wff with all these operators deleted except the last.

Considerations of space preclude an account of the many other axiomatic systems of modal logic that have been investigated. Some of these are weaker than **T**; such systems normally contain the axioms of **T** either as axioms or as theorems but have only a restricted form of the rule of necessitation. Another group comprises systems that are stronger than **S4** but weaker than **S5**; some of these have proved fruitful in developing a logic of temporal relations. Yet another group includes systems that are stronger than **S4** but independent of **S5** in the sense explained above.

Modal predicate logics can be formed also by making analogous additions to LPC instead of to PC.

Validity in modal logic. The task of defining validity for modal wffs is complicated by the fact that, even if the truth values of all of the variables in a wff are given, it is not obvious how one should set about calculating the truth value of the whole wff. Nevertheless, a number of defini-

The system **T**

tions of validity applicable to modal wffs have been given, each of which turns out to match some axiomatic modal system, in the sense that it brings out as valid those wffs, and no others, that are theorems of that system. Most, if not all, of these accounts of validity can be thought of as variant ways of giving formal precision to the idea that necessity is truth in every "possible world" or "conceivable state of affairs." The simplest such definition is this: Let a model be constructed by first assuming a (finite or infinite) set W of "worlds." In each world, independently of all the others, let each propositional variable then be assigned either the value 1 or the value 0. In each world the values of truth functions are calculated in the usual way from the values of their arguments in that world. In each world, however, $L\alpha$ is to have the value 1 if α has the value 1 not only in that world but in every other world in W as well and is otherwise to have the value 0; and in each world $M\alpha$ is to have the value 1 if α has value 1 either in that world or in some other world in W and is otherwise to have the value 0. These rules enable one to calculate a value (1 or 0) in any world in W for any given wff, once the values of the variables in each world in W are specified. A model is defined as consisting of a set of worlds together with a value assignment of the kind just described. A wff is valid if and only if it has the value 1 in every world in every model. It can be proved that the wffs that are valid by this criterion are precisely the theorems of **S5**; for this reason models of the kind here described may be called **S5**-models, and validity as just defined may be called **S5**-validity.

A definition of **T**-validity (*i.e.*, one that can be proved to bring out as valid precisely the theorems of **T**) can be given as follows: a **T**-model consists of a set of worlds W and a value assignment to each variable in each world, as before. It also includes a specification, for each world in W , of some subset of W as the worlds that are "accessible" to that world. Truth functions are evaluated as before; but in each world in the model, $L\alpha$ is to have the value 1 if α has value 1 in that world and in every other world in W accessible to it and is otherwise to have value 0. And in each world, $M\alpha$ is to have the value 1 if α has value 1 either in that world or in some other world accessible to it and is otherwise to have value 0. (In other words, in computing the value of $L\alpha$ or $M\alpha$ in a given world, no account is taken of the value of α in any other world not accessible to it.) A wff is **T**-valid if and only if it has the value 1 in every world in every **T**-model.

An **S4**-model is defined as is a **T**-model except that it is required that the accessibility relation be transitive—*i.e.*, that, where $w_1, w_2,$ and w_3 are any worlds in W , if w_1 is accessible to w_2 and w_2 is accessible to w_3 , then w_1 is accessible to w_3 . A wff is **S4**-valid if and only if it has the value 1 in every world in every **S4**-model. The **S4**-valid wffs can be shown to be precisely the theorems of **S4**. Finally, a definition of validity is obtained that will match the system **B** by requiring that the accessibility relation be symmetrical but not that it be transitive.

For all four systems, effective decision procedures for validity can be given. Further modifications of the general method described have yielded validity definitions that match many other axiomatic modal systems, and the method can be adapted to give a definition of validity for intuitionistic **PC**: For a number of axiomatic modal systems, however, no satisfactory account of validity has been devised. Validity can also be defined for various modal predicate logics by combining the definition of **LPC**-validity given above in *Validity in LPC* with the relevant accounts of validity for modal systems, but a modal logic based on **LPC** is, like **LPC** itself, an undecidable system.

SET THEORY

Only a sketchy account of set theory is given here; for further information see the article **SET THEORY**. Set theory is a logic of classes—*i.e.*, of collections (finite or infinite) or aggregations of objects of any kind, which are known as the members of the classes in question. Some logicians use the terms "class" and "set" interchangeably; others distinguish between them, defining a set (for example) as a class that is itself a member of some class and defining a

proper class as one that is not a member of any class. It is usual to write " \in " for "is a member of" and to abbreviate $\neg(x \in y)$ to $x \notin y$. A particular class may be specified either by listing all its members or by stating some condition of membership, in which (latter) case the class consists of all those things and only those that satisfy that condition (it is used, for example, when one speaks of the class of inhabitants of London or the class of prime numbers). Clearly, the former method is available only for finite classes and may be very inconvenient even then; the latter, however, is of more general applicability. Two classes that have precisely the same members are regarded as the same class or are said to be identical with each other, even if they are specified by different conditions: *i.e.*, identity of classes is identity of membership, not identity of specifying conditions. This principle is known as the principle of extensionality. A class with no members, such as the class of Chinese popes, is said to be null. Since the membership of all such classes is the same, there is only one null class, which is therefore usually called *the* null class (or sometimes the empty class); it is symbolized by Λ or \emptyset . The notation $x = y$ is used for " x is identical with y ," and $\neg(x = y)$ is usually abbreviated to $x \neq y$. The expression $x = \Lambda$ therefore means that the class x has no members, and $x \neq \Lambda$ means that x has at least one member.

A member of a class may itself be a class. The class of dogs, for example, is a member of the class of species of animals. An individual dog, however, though a member of the former class, is not a member of the latter—because an individual dog is not a species of animal (if the number of dogs increases, the number of species of animals does not thereby increase). Class membership is therefore not a transitive relation. The relation of class inclusion, however (to be carefully distinguished from class membership), is transitive. A class x is said to be included in a class y (written: $x \subseteq y$) if and only if every member of x is also a member of y . (This is not meant to exclude the possibility that x and y may be identical.) If x is included in, but is not identical with, y —*i.e.*, if every member of x is a member of y but some members of y are not members of x — x is said to be properly included in y (written: $x \subset y$).

It is perhaps natural to assume that for every storable condition there is a class (null or otherwise) of objects that satisfy that condition. This assumption is known as the principle of comprehension. In the unrestricted form just mentioned, however, this principle has been found to lead to inconsistencies and hence cannot be accepted as it stands. One storable condition, for example, is non-self-membership—*i.e.*, the property possessed by a class if and only if it is not a member of itself. This in fact appears to be a condition that most classes do fulfill; the class of dogs, for example, is not itself a dog and hence is not a member of the class of dogs.

Let it now be assumed that the class of all classes that are not members of themselves can be formed and let this class be y . Then any class x will be a member of y if and only if it is not a member of itself; *i.e.*, for any class x , $(x \in y) \equiv (x \notin x)$. The question can then be asked whether y is a member of itself or not, with the following awkward result: If it is a member of itself, then it fails to fulfill the condition of membership of y , and hence it is not a member of y —*i.e.*, not a member of itself. On the other hand, if y is not a member of itself, then it does fulfill the required condition, and therefore it is a member of y —*i.e.*, of itself. Hence the equivalence $(y \in y) \equiv (y \notin y)$ results, which is self-contradictory. This perplexing conclusion, which was pointed out by Bertrand Russell, is known as Russell's paradox. Russell's own solution to it and to other similar difficulties was to regard classes as forming a hierarchy of types and to posit that a class could only be regarded sensibly as a member, or a nonmember, of a class at the next higher level in the hierarchy. The effect of this theory is to make $x \in x$, and therefore $x \notin x$, ill-formed. Another kind of solution, however, is based upon the distinction made earlier between two kinds of classes, those that are sets and those that are not—a set being defined as a class that is itself a member of some class. The unrestricted principle of comprehension is then replaced by the weaker principle that for every condition

Models for the four systems

Adaptations to other systems

Interrelations of classes

Russell's paradox: theory of types

there is a class the members of which are the individuals or sets that fulfill that condition. Other solutions have also been devised, but none has won universal acceptance, with the result that several different versions of set theory are found in the literature of the subject.

Formally, set theory can be derived by the addition of various special axioms to a rather modest form of LPC that contains no predicate variables and only a single primitive dyadic predicate constant (\in) to represent membership. Sometimes LPC-with-identity is used, and there are then two primitive dyadic predicate constants (\in and $=$). In some versions the variables x, y, \dots are taken to range only over sets or classes; in other versions they range over individuals as well. The special axioms vary, but the basis normally includes the principle of extensionality and some restricted form of the principle of comprehension, or some elements from which these can be deduced.

Notation, definitions, theorems

A notation to express theorems about classes can be either defined in various ways (not detailed here) in terms of the primitives mentioned above or else introduced independently. The main elements of one widely used notation are the following: if α is an expression containing some free occurrence of x , the expression $\{x : \alpha\}$ is used to stand for the class of objects fulfilling the condition expressed by α —e.g., $\{x : x \text{ is a prime number}\}$ represents the class of prime numbers; $\{x\}$ represents the class the only member of which is x ; $\{x, y\}$ the class the only members of which are x and y ; and so on. $\langle x, y \rangle$ represents the class the members of which are x and y in that order (thus $\{x, y\}$ and $\{y, x\}$ are identical; but $\langle x, y \rangle$ and $\langle y, x \rangle$ are in general not identical). Let x and y be any classes, as (for example) those of the dots on the two arms of a stippled cross. The intersection of x and y , symbolized as $x \cap y$, is the class the members of which are the objects common to x and y —in this case the dots within the area where the arms cross—i.e., $\{z : z \in x \cdot z \in y\}$. Similarly, the union of x and y , symbolized as $x \cup y$, is the class the members of which are the members of x together with those of y —in this case all the dots on the cross—i.e., $\{z : z \in x \vee z \in y\}$; the complement of x , symbolized as \bar{x} , is the class the members of which are all those objects that are not members of x —i.e., $\{y : y \notin x\}$; the complement of y in x , symbolized as $x - y$, is the class of all objects that are members of x but not of y —i.e., $\{z : z \in x \cdot z \notin y\}$; the universal class, symbolized as V , is the class of which everything is a member, definable as the complement of the null class—i.e., as $\bar{\Lambda}$. Λ itself is sometimes taken as a primitive individual constant, sometimes defined as $\{x : x \neq x\}$ —the class of objects that are not identical with themselves.

Among the simpler theorems of set theory are

$$\begin{aligned} &(\forall x)(x \cap x = x), \\ &(\forall x)(\forall y)(x \cap y = y \cap x); \end{aligned}$$

and corresponding theorems for \cup :

$$\begin{aligned} &(\forall x)(\forall y)(\forall z)[x \cap (y \cup z) = (x \cap y) \cup (x \cap z)], \\ &(\forall x)(\forall y)[\bar{(x \cap y)} = \bar{x} \cup \bar{y}]; \end{aligned}$$

and corresponding theorems with \cap and \cup interchanged:

$$\begin{aligned} &(\forall x)(\bar{\bar{x}} = x), \\ &(\forall x)(\forall y)(x - y = x \cap \bar{y}), \\ &(\forall x)(\Lambda \subset x), \\ &(\forall x)(x \cap \Lambda = \Lambda), \\ &(\forall x)(x \cup \Lambda = x). \end{aligned}$$

In these theorems, the variables range over classes. In several cases, there are obvious analogies to valid wffs of PC.

Apart from its own intrinsic interest, set theory has an importance for the foundations of mathematics in that it is widely held that the natural numbers can be adequately defined in set-theoretic terms. Moreover, given suitable axioms, standard postulates for natural-number arithmetic can be derived as theorems within set theory.

(G.E.H./M.L.Sc.)

Metalogic

Metalogic is the study of the syntax and the semantics of formal languages and formal systems. It is related to, but does not include, the formal treatment of natural lan-

guages. (For a discussion of the syntax and semantics of natural languages, see the article LINGUISTICS.)

NATURE, ORIGINS, AND INFLUENCES OF METALOGIC

Syntax and semantics. A formal language usually requires a set of formation rules—i.e., a complete specification of the kinds of expressions that shall count as well-formed formulas (sentences or meaningful expressions), applicable mechanically, in the sense that a machine could check whether a candidate satisfies the requirements. This specification usually contains three parts: (1) a list of primitive symbols (basic units) given mechanically, (2) certain combinations of these symbols, singled out mechanically as forming the simple (atomic) sentences, and (3) a set of inductive clauses—inductive inasmuch as they stipulate that natural combinations of given sentences formed by such logical connectives as the disjunction “or,” which is symbolized “ \vee ”; “not,” symbolized “ \sim ”; and “for all ___,” symbolized “ $(\forall _)$,” are again sentences. [“($\forall _$)” is called a quantifier, as is also “there is some ___,” symbolized “ $(\exists _)$.”] Since these specifications are concerned only with symbols and their combinations and not with meanings, they involve only the syntax of the language.

Specification of a formal language

An interpretation of a formal language is determined by formulating an interpretation of the atomic sentences of the language with regard to a domain of objects—i.e., by stipulating which objects of the domain are denoted by which constants of the language and which relations and functions are denoted by which predicate letters and function symbols. The truth-value (whether “true” or “false”) of every sentence is thus determined according to the standard interpretation of logical connectives. For example, $p \cdot q$ is true if and only if p and q are true. (Here, the dot means the conjunction “and,” not the multiplication operation “times.”) Thus, given any interpretation of a formal language, a formal concept of truth is obtained. Truth, meaning, and denotation are semantic concepts.

If, in addition, a formal system in a formal language is introduced, certain syntactic concepts arise—namely, axioms, rules of inference, and theorems. Certain sentences are singled out as axioms. These are (the basic) theorems. Each rule of inference is an inductive clause, stating that, if certain sentences are theorems, then another sentence related to them in a suitable way is also a theorem. If p and “either not- p or q ” ($\bar{p} \vee q$) are theorems, for example, then q is a theorem. In general, a theorem is either an axiom or the conclusion of a rule of inference whose premises are theorems.

In 1931 Gödel made the fundamental discovery that, in most of the interesting (or significant) formal systems, not all true sentences are theorems. It follows from this finding that semantics cannot be reduced to syntax: thus syntax, which is closely related to proof theory, must often be distinguished from semantics, which is closely related to model theory. Roughly speaking, syntax—as conceived in the philosophy of mathematics—is a branch of number theory, and semantics is a branch of set theory, which deals with the nature and relations of aggregates.

Historically, as logic and axiomatic systems became more and more exact, there emerged, in response to a desire for greater lucidity, a tendency to pay greater attention to the syntactic features of the languages employed rather than to concentrate exclusively on intuitive meanings. In this way, logic, the axiomatic method (such as that employed in geometry), and semiotic (the general science of signs) converged toward metalogic.

The axiomatic method. The best known axiomatic system is that of Euclid for geometry. In a manner similar to that of Euclid, every scientific theory involves a body of meaningful concepts and a collection of true or believed assertions. The meaning of a concept can often be explained or defined in terms of other concepts, and, similarly, the truth of an assertion or the reason for believing it can usually be clarified by indicating that it can be deduced from certain other assertions already accepted. The axiomatic method proceeds in a sequence of steps, beginning with a set of primitive concepts and propositions and then defining or deducing all other concepts and propositions in the theory from them.

The realization which arose in the 19th century that there are different possible geometries led to a desire to separate abstract mathematics from spatial intuition; in consequence, many hidden axioms were uncovered in Euclid's geometry. These discoveries were organized into a more rigorous axiomatic system by David Hilbert in his *Grundlagen der Geometrie* (1899; *The Foundations of Geometry*). In this and related systems, however, logical connectives and their properties are taken for granted and remain implicit. If the logic involved is taken to be that of the predicate calculus, the logician can then arrive at such formal systems as that discussed above.

Once such formal systems are obtained, it is possible to transform certain semantic problems into sharper syntactic problems. It has been asserted, for example, that non-Euclidean geometries must be self-consistent systems because they have models (or interpretations) in Euclidean geometry, which in turn has a model in the theory of real numbers. It may then be asked, however, how it is known that the theory of real numbers is consistent in the sense that no contradiction can be derived within it. Obviously, modeling can establish only a relative consistency and has to come to a stop somewhere. Having arrived at a formal system (say, of real numbers), however, the consistency problem then has the sharper focus of a syntactic problem: that of considering all the possible proofs (as syntactic objects) and asking whether any of them ever has (say) $0 = 1$ as the last sentence.

As another example, the question whether a system is categorical—that is, whether it determines essentially a unique interpretation in the sense that any two interpretations are isomorphic—may be explored. This semantic question can to some extent be replaced by a related syntactic question, that of completeness: whether there is in the system any sentence having a definite truth-value in the intended interpretation such that neither that sentence nor its negation is a theorem. Even though it is now known that the semantic and syntactic concepts are different, the vague requirement that a system be “adequate” is clarified by both concepts. The study of such sharp syntactic questions as those of consistency and completeness, which was emphasized by Hilbert, was named “metamathematics” (or “proof theory”) by him about 1920.

Logic and metalogic. In one sense, logic is to be identified with the predicate calculus of the first order, the calculus in which the variables are confined to individuals of a fixed domain—though it may include as well the logic of identity, symbolized “=,” which takes the ordinary properties of identity as part of logic. In this sense Gottlob Frege achieved a formal calculus of logic as early as 1879. Sometimes logic is construed, however, as including also higher-order predicate calculi, which admit variables of higher types, such as those ranging over predicates (or classes and relations) and so on. But then it is a small step to the inclusion of set theory, and, in fact, axiomatic set theory is often regarded as a part of logic. For the purposes of this article, however, it is more appropriate to confine the discussion to logic in the first sense.

It is hard to separate out significant findings in logic from those in metalogic, because all theorems of interest to logicians are about logic and therefore belong to metalogic. If p is a mathematical theorem—in particular, one about logic—and P is the conjunction of the mathematical axioms employed for proving p , then every p can be turned into a theorem, “either not- P or p ,” in logic. Mathematics is not done, however, by carrying out explicitly all the steps as formalized in logic; the selection and intuitive grasp of the axioms is important both for mathematics and for metamathematics. Actual derivations in logic, such as those carried out just prior to World War I by Alfred North Whitehead and Bertrand Russell, are of little intrinsic interest to logicians. It might therefore appear redundant to introduce the term metalogic. In the present classification, however, metalogic is conceived as dealing not only with findings about logical calculi but also with studies of formal systems and formal languages in general.

An ordinary formal system differs from a logical calculus in that the system usually has an intended interpretation, whereas the logical calculus deliberately leaves the possible

interpretations open. Thus, one speaks, for example, of the truth or falsity of sentences in a formal system, but with respect to a logical calculus one speaks of validity (*i.e.*, being true in all interpretations or in all possible worlds) and of satisfiability (or having a model—*i.e.*, being true in some particular interpretation). Hence, the completeness of a logical calculus has quite a different meaning from that of a formal system: a logical calculus permits many sentences such that neither the sentence nor its negation is a theorem because it is true in some interpretations and false in others, and it requires only that every valid sentence be a theorem.

Semiotic. Originally, the word “semiotic” meant the medical theory of symptoms; however, an empiricist, John Locke, used the term in the 17th century for a science of signs and significations. The current usage was recommended especially by Rudolf Carnap—see his *Introduction to Semantics* (1942) and his reference there to Charles William Morris, who suggested a threefold distinction. According to this usage, semiotic is the general science of signs and languages, consisting of three parts: (1) pragmatics (in which reference is made to the user of the language), (2) semantics (in which one abstracts from the user and analyzes only the expressions and their meanings), and (3) syntax (in which one abstracts also from the meanings and studies only the relations between expressions).

Considerable effort since the 1970s has gone into the attempt to formalize some of the pragmatics of natural languages. The use of indexical expressions to incorporate reference to the speaker, his or her location, or the time of either the utterance or the events mentioned was of little importance to earlier logicians, who were primarily interested in universal truths or mathematics. With the increased interest in linguistics there has come an increased effort to formalize pragmatics.

At first Carnap exclusively emphasized syntax. But gradually he came to realize the importance of semantics, and the door was thus reopened to many difficult philosophical problems.

Certain aspects of metalogic have been instrumental in the development of the approach to philosophy commonly associated with the label of logical positivism. In his *Tractatus Logico-Philosophicus* (1922; originally published under another title, 1921), Ludwig Wittgenstein, a seminal thinker in the philosophy of language, presented an exposition of logical truths as sentences that are true in all possible worlds. One may say, for example, “It is raining or it is not raining,” and in every possible world one of the disjuncts is true. On the basis of this observation and certain broader developments in logic, Carnap tried to develop formal treatments of science and philosophy.

It has been thought that the success that metalogic had achieved in the mathematical disciplines could be carried over into physics and even into biology or psychology. In so doing, the logician gives a branch of science a formal language in which there are logically true sentences having universal logical ranges and factually true sentences having universal logical ranges and factually true ones having more restricted ranges. (Roughly speaking, the logical range of a sentence is the set of all possible worlds in which it is true.)

A formal solution of the problem of meaning has also been proposed for these disciplines. Given the formal language of a science, it is possible to define a notion of truth. Such a truth definition determines the truth condition for every sentence—*i.e.*, the necessary and sufficient conditions for its truth. The meaning of a sentence is then identified with its truth condition because, as Carnap wrote:

To understand a sentence, to know what is asserted by it, is the same as to know under what conditions it would be true. . . . To know the truth condition of a sentence is (in most cases) much less than to know its truth-value, but it is the necessary starting point for finding out its truth-value.

Influences in other directions. Metalogic has led to a great deal of work of a mathematical nature in axiomatic set theory, model theory, and recursion theory (in which functions that are computable in a finite number of steps are studied).

In a different direction, the devising of Turing computing

Trans-
formations
from
semantics
to syntax

Formal
systems
and logical
calculi

Logical
positivism
and the
formaliza-
tion of
science

Turing machines, philosophy of logic, and ontology

machines, involving abstract designs for the explication of mechanical logical procedures, has led to the investigation of idealized computers, with ramifications in the theory of finite automata and mathematical linguistics.

Among philosophers of language, there is a widespread tendency to stress the philosophy of logic. The contrast, for example, between intensional concepts and extensional concepts; the role of meaning in natural languages as providing truth conditions; the relation between formal and natural logic (*i.e.*, the logic of natural languages); and the relation of ontology, the study of the kinds of entities that exist, to the use of quantifiers—all these areas are receiving extensive consideration and discussion. There are also efforts to produce formal systems for empirical sciences such as physics, biology, and even psychology. Many scholars have doubted, however, whether these latter efforts have been fruitful.

NATURE OF A FORMAL SYSTEM AND OF ITS FORMAL LANGUAGE

Example of a formal system. In order to clarify the abstract concepts of metalogic, a formal system *N* (with its formal language) may be considered for illustration.

Formation rules. The system may be set up by employing the following formation rules:

1. The following are primitive symbols: “~,” “∨,” “∀,” and “=” and the symbols used for grouping, “(” and “)”; the function symbols for “successor,” “*S*,” and for arithmetical addition and multiplication, “+” and “·”; constants 0, 1; and variables *x*, *y*, *z*,
2. The following are terms: a constant is a term; a variable is a term; if *a* is a term, *Sa* is a term; and, if *a* and *b* are terms, *a* + *b* and *a* · *b* are terms.
3. Atomic sentences are thus specified: if *a* and *b* are terms, *a* = *b* is a sentence.
4. Other sentences can be defined as follows: if *A* and *B* are sentences and *v* is a variable, then ~*A*, *A* ∨ *B*, and (∀*v*)*A* are sentences.

Axioms and rules of inference. The system may be developed by adopting certain sentences as axioms and following certain rules of inference.

1. The basic axioms and rules are to be those of the first-order predicate calculus with identity (see above *Formal logic*).
2. The following additional axioms of *N* are stipulated:
 - a. Zero (0) is not a successor:

$$\sim Sx = 0$$

- b. No two different numbers have the same successor:

$$\sim (Sx = Sy) \vee x = y$$

- c. Recursive definition of addition:

$$\begin{aligned} x + 0 &= x \\ x + Sy &= S(x + y) \end{aligned}$$

(From this, with the understanding that 1 is the successor of 0, one can easily show that *Sx* = *x* + 1.)

- d. Recursive definition of multiplication:

$$\begin{aligned} x \cdot 0 &= 0 \\ x \cdot Sy &= (x \cdot y) + x \end{aligned}$$

3. Rule of inference (the principle of mathematical induction): If zero has some property *p* and it is the case that if any number has *p* then its successor does, then every number has *p*. With some of the notation from above, this can be expressed: If *A*(0) and (∀*x*)(~*A*(*x*) ∨ *A*(*Sx*)) are theorems, then (∀*x*)*A*(*x*) is a theorem.

The system *N* as specified by the foregoing rules and axioms is a formal system in the sense that, given any combination of the primitive symbols, it is possible to check mechanically whether it is a sentence of *N*, and, given a finite sequence of sentences, it is possible to check mechanically whether it is a (correct) proof in *N*—*i.e.*, whether each sentence either is an axiom or follows from preceding sentences in the sequence by a rule of inference. Viewed in this way, a sentence is a theorem if and only if there exists a proof in which it appears as the last sentence. It is not required of a formal system, however, that it be

possible to decide mechanically whether or not a given sentence is a theorem; and, in fact, it has been proved that no such mechanical method exists.

Truth definition of the given language. The formal system *N* admits of different interpretations, according to findings of Gödel (from 1931) and of the Norwegian mathematician Thoralf Skolem, a pioneer in metalogic (from 1933). The originally intended, or standard, interpretation takes the ordinary nonnegative integers {0, 1, 2, . . .} as the domain, the symbols 0 and 1 as denoting zero and one, and the symbols + and · as standing for ordinary addition and multiplication. Relative to this interpretation, it is possible to give a truth definition of the language of *N*.

It is necessary first to distinguish between open and closed sentences. An open sentence, such as *x* = 1, is one that may be either true or false depending on the value of *x*, but a closed sentence, such as 0 = 1 and (∀*x*)(*x* = 0) or “All *x*’s are zero,” is one that has a definite truth-value—in this case, false (in the intended interpretation).

Open and closed sentences

1. A closed atomic sentence is true if and only if it is true in the intuitive sense; for example, 0 = 0 is true, 0 + 1 = 0 is false.

This specification as it stands is not syntactic, but, with some care, it is possible to give an explicit and mechanical specification of those closed atomic sentences that are true in the intuitive sense.

2. A closed sentence ~*A* is true if and only if *A* is not true.
3. A closed sentence *A* ∨ *B* is true if and only if either *A* or *B* is true.
4. A closed sentence (∀*v*)*A*(*v*) is true if and only if *A*(*v*) is true for every value of *v*—*i.e.*, if *A*(0), *A*(1), *A*(1 + 1), . . . are all true.

The above definition of truth is not an explicit definition; it is an inductive one. Using concepts from set theory, however, it is possible to obtain an explicit definition that yields a set of sentences that consists of all the true ones and only them. If Gödel’s method of representing symbols and sentences by numbers is employed, it is then possible to obtain in set theory a set of natural numbers that are just the Gödel numbers of the true sentences of *N*.

There is a definite sense in which it is impossible to define the concept of truth within a language itself. This is proved by the liar paradox: if the sentence “I am lying,” or alternatively

$$\text{This sentence is not true.} \tag{1}$$

is considered, it is clear—since (1) is “This sentence”—that if (1) is true, then (1) is false; on the other hand, if (1) is false, then (1) is true. In the case of the system *N*, if the concept of truth were definable in the system itself, then (using a device invented by Gödel) it would be possible to obtain in *N* a sentence that amounts to (1) and that thereby yields a contradiction.

DISCOVERIES ABOUT FORMAL MATHEMATICAL SYSTEMS

The two central questions of metalogic are those of the completeness and consistency of a formal system based on axioms. In 1931 Gödel made fundamental discoveries in these areas for the most interesting formal systems. In particular, he discovered that, if such a system is ω-consistent—*i.e.*, devoid of contradiction in a sense to be explained below—then it is not complete and that, if a system is consistent, then the statement of its consistency, easily expressible in the system, is not provable in it.

Soon afterward, in 1934, Gödel modified a suggestion that had been offered by Jacques Herbrand, a French mathematician, and introduced a general concept of recursive functions—*i.e.*, of functions mechanically computable by a finite series of purely combinatorial steps. In 1936 Alonzo Church, a mathematical logician, Alan Mathison Turing, originator of a theory of computability, and Emil L. Post, a specialist in recursive unsolvability, all argued for this concept (and certain equivalent notions), thereby arriving at stable and exact conceptions of “mechanical,” “computable,” “recursive,” and “formal” that explicate the intuitive concept of what a mechanical computing procedure is. As a result of the development of recursion theory, it is now possible to prove not only that certain classes of problems are mechanically solvable

Recursive functions and decidability

(which could be done without the theory) but also that certain others are mechanically unsolvable (or absolutely unsolvable). The most notable example of such unsolvability is the discovery, made in 1970, that there is no algorithm, or rule of repetitive procedure, for determining which Diophantine equations (*i.e.*, polynomial equations of which the coefficients are whole numbers) have integer solutions. This solution gives a negative solution to the 10th problem in the famous list presented by Hilbert at the International Mathematical Congress in 1900.

In this way, logicians have finally arrived at a sharp concept of a formal axiomatic system, because it is no longer necessary to leave "mechanical" as a vague non-mathematical concept. In this way, too, they have arrived at sharp concepts of decidability. In one sense, decidability is a property of sets (of sentences): that of being subject (or not) to mechanical methods by which to decide in a finite number of steps, for any closed sentence of a given formal system (*e.g.*, of \mathbb{N}), whether it is true or not or—as a different question—whether it is a theorem or not. In another sense, decidability can refer to a single closed sentence: the sentence is called undecidable in a formal system if neither it nor its negation is a theorem. Using this concept, Gödel's incompleteness theorem is sometimes stated thus: "Every interesting (or significant) formal system has some undecidable sentences."

Given these developments, it was easy to extend Gödel's findings, as Church did in 1936, to show that interesting formal systems such as \mathbb{N} are undecidable (both with regard to theorems and with regard to true sentences).

The two incompleteness theorems. The first and most central finding in this field is that systems such as \mathbb{N} are incomplete and incompletable because Gödel's theorem applies to any reasonable and moderately rich system. The proof of this incompleteness may be viewed as a modification of the liar paradox, which shows that truth cannot be defined in the language itself. Since provability in a formal system can often be expressed in the system itself, one is led to the conclusion of incompleteness.

Let us consider the sentence

This sentence is not provable in the system. (2)

In particular, \mathbb{N} may be thought of as the system being studied. Representing expressions by numbers and using an ingenious substitution function, Gödel was able to find in the system a sentence p that could be viewed as expressing (2).

Once such a sentence is obtained, some strong conclusions result. If the system is complete, then either the sentence p or its negation is a theorem of the system. If p is a theorem, then intuitively p or (2) is false, and there is in some sense a false theorem in the system. Similarly, if $\sim p$ is a theorem, then it says that $\sim(2)$ or that p is provable in the system. Since $\sim p$ is a theorem, it should be true, and there seem then to be two conflicting sentences that are both true—namely, p is provable in the system and $\sim p$ is provable in it. This can be the case only if the system is inconsistent.

A careful examination of this inexact line of reasoning leads to Gödel's exact theorem, which says that, if a system is reasonably rich and ω -consistent, then p is undecidable in it. The notion of ω -consistency is stronger than consistency, but it is a very reasonable requirement, since it demands merely that one cannot prove in a system both that some number does not have the property A and yet for each number that it does have the property A —*i.e.*, that $(\exists x)\sim A(x)$ and also all of $A(0), A(1), \dots$ are theorems. The American mathematician J. Barkley Rosser, who also contributed to number theory and applied mathematics, weakened the hypothesis to mere consistency in 1936, at the expense of complicating somewhat the initial sentence (2).

More exactly, Gödel showed that, if the system is consistent, then p is not provable; if it is ω -consistent, then $\sim p$ is not provable. The first half leads to Gödel's theorem on consistency proofs, which says that if a system is consistent, then the arithmetic sentence expressing the consistency of the system cannot be proved in the system. This is usually stated briefly thus: that no interesting sys-

tem can prove its own consistency or that there exists no consistency proof of a system that can be formalized in the system itself.

The proof of this theorem consists essentially of a formalization in arithmetic of the arithmetized version of the proof of the statement, "If a system is consistent, then p is not provable"; *i.e.*, it consists of a derivation within number theory of p itself from the arithmetic sentence that says that the system is consistent. Hence, if the arithmetic sentence were provable, p would also be provable—contradicting the previous result. This proof, which was only briefly outlined by Gödel, has been carried out in detail by Paul Bernays in his joint work with Hilbert. Moreover, the undecidable sentence p is always of a relatively simple form—namely, the form $(\forall x)A(x)$, "For every x , x is A ," in which A is a recursive, in fact a primitive recursive, predicate.

Decidability and undecidability. The first incompleteness theorem yields directly the fact that truth in a system (*e.g.*, in \mathbb{N}) to which the theorem applies is undecidable. If it were decidable, then all true sentences would form a recursive set, and they could be taken as the axioms of a formal system that would be complete. This claim depends on the reasonable and widely accepted assumption that all that is required of the axioms of a formal system is that they make it possible to decide effectively whether a given sentence is an axiom.

Alternatively, the above assumption can be avoided by resorting to a familiar lemma, or auxiliary truth: that all recursive or computable functions and relations are representable in the system (*e.g.*, in \mathbb{N}). Since truth in the language of a system is itself not representable (definable) in the system, it cannot, by the lemma, be recursive (*i.e.*, decidable).

The same lemma also yields the undecidability of such systems with regard to theorems. Thus, if there were a decision procedure, there would be a computable function f such that $f(i)$ equals 1 or 0 according as the i th sentence is a theorem or not. But then what $f(i) = 0$ says is just that the i th sentence is not provable. Hence, using Gödel's device, a sentence (say the t th) is again obtained saying of itself that it is not provable. If $f(t) = 0$ is true, then, because f is representable in the system, it is a theorem of the system. But then, because $f(t) = 0$ is (equivalent to) the t th sentence, $f(t) = 1$ is also true and therefore provable in the system. Hence, the system, if consistent, is undecidable with regard to theorems.

Although the system \mathbb{N} is incompletable and undecidable, it has been discovered by the Polish logician M. Presburger and by Skolem (both in 1930) that arithmetic with addition alone or multiplication alone is decidable (with regard to truth) and therefore has complete formal systems. Another well-known positive finding is that of the Polish-American semanticist and logician Alfred Tarski, who developed a decision procedure for elementary geometry and elementary algebra (1951).

Consistency proofs. The best-known consistency proof is that of the German mathematician Gerhard Gentzen (1936) for the system \mathbb{N} of classical (or ordinary, in contrast to intuitionistic) number theory. Taking ω (omega) to represent the next number beyond the natural numbers (called the "first transfinite number"), Gentzen's proof employs an induction in the realm of transfinite numbers ($\omega + 1, \omega + 2, \dots; 2\omega, 2\omega + 1, \dots; \omega^2, \omega^2 + 1, \dots$), which is extended to the first epsilon-number, ϵ_0 (defined as the limit of $\omega, \omega^\omega, \omega^{\omega^\omega}, \dots$), which is not formalizable in \mathbb{N} . This proof, which has appeared in several variants, has opened up an area of rather extensive work.

Intuitionistic number theory, which denies the classical concept of truth and consequently eschews certain general laws such as "either A or $\sim A$," and its relation to classical number theory have also been investigated (see MATHEMATICS, THE FOUNDATIONS OF: *Intuitionism*). This investigation is considered significant, because intuitionism is believed to be more constructive and more evident than classical number theory. In 1932 Gödel found an interpretation of classical number theory in the intuitionistic theory (also found by Gentzen and by Bernays). In 1958 Gödel extended his findings to obtain constructive

Undecidability of truth and theoremhood

Gödel's exact theorem

interpretations of sentences of classical number theory in terms of primitive recursive functionals.

More recently, work has been done to extend Gentzen's findings to ramified theories of types and to fragments of classical analysis and to extend Gödel's interpretation and to relate classical analysis to intuitionistic analysis. Also, in connection with these consistency proofs, various proposals have been made to present constructive notations for the ordinals of segments of the German mathematician Georg Cantor's second number class, which includes ω and the first epsilon-number and much more. A good deal of discussion has been devoted to the significance of the consistency proofs and the relative interpretations for epistemology (the theory of knowledge).

DISCOVERIES ABOUT LOGICAL CALCULI

The calculi of formal logic. The two main branches of formal logic are the propositional calculus and the predicate calculus (see above *Formal logic*).

The propositional calculus. It is easy to show that the propositional calculus is complete in the sense that every valid sentence in it—*i.e.*, every tautology, or sentence true in all possible worlds (in all interpretations)—is a theorem, as may be seen in the following example. "Either p or not- p " ($p \vee \sim p$) is always true because p is either true or false. In the former case, $p \vee \sim p$ is true because p is true; in the latter case, because $\sim p$ is true. One way to prove the completeness of this calculus is to observe that it is sufficient to reduce every sentence to a conjunctive normal form—*i.e.*, to a conjunction of disjunctions of single letters and their negations. But any such conjunction is valid if and only if every conjunct is valid; and a conjunct is valid if and only if it contains some letter p as well as $\sim p$ as parts of the whole disjunction. Completeness follows because (1) such conjuncts can all be proved in the calculus and (2) if these conjuncts are theorems, then the whole conjunction is also a theorem.

The consistency of the propositional calculus (its freedom from contradiction) is more or less obvious, because it can easily be checked that all its axioms are valid—*i.e.*, true in all possible worlds—and that the rules of inference carry from valid sentences to valid sentences. But a contradiction is not valid; hence, the calculus is consistent. The conclusion, in fact, asserts more than consistency, for it holds that only valid sentences are provable.

The calculus is also easily decidable. Since all valid sentences, and only these, are theorems, each sentence can be tested mechanically by putting true and false for each letter in the sentence. If there are n letters, there are 2^n possible substitutions. A sentence is then a theorem if and only if it comes out true in every one of the 2^n possibilities.

The independence of the axioms is usually proved by using more than two truth values. These values are divided into two classes: the desired and the undesired. The axiom to be shown independent can then acquire some undesired value, whereas all the theorems that are provable without this axiom always get the desired values. This technique is what originally suggested the many-valued logics.

The first-order predicate calculus. The problem of consistency for the predicate calculus is relatively simple. A world may be assumed in which there is only one object a . In this case, both the universally quantified and the existentially quantified sentences ($\forall x)A(x)$ and ($\exists x)A(x)$ reduce to the simple sentence $A(a)$, and all quantifiers can be eliminated. It may easily be confirmed that, after the reduction, all theorems of the calculus become tautologies (*i.e.*, theorems in the propositional calculus). If F is any predicate, such a sentence as "Every x is F and not every x is F "—*i.e.*, ($\forall x)F(x) \cdot \sim(\forall x)F(x)$ —is then reduced to " a is both A and not- A "— $A(a) \cdot \sim A(a)$ —which is not a tautology; therefore, the original sentence is not a theorem; hence, no contradiction can be a theorem. If F is simple, then F and A are the same. If F is complex and contains ($\forall y$) or ($\exists z$), etc., then A is the result obtained by iterating the transformation of eliminating ($\forall y$), etc. In fact, it can be proved quite directly not only that the calculus is consistent but also that all its theorems are valid.

The discoveries that the calculus is complete and undecidable are much more profound than the discovery of

its consistency. Its completeness was proved by Gödel in 1930; its undecidability was established with quite different methods by Church and Turing in 1936. Given the general developments that occurred up to 1936, its undecidability also follows in another way from Theorem X of Gödel's paper of 1931.

Completeness means that every valid sentence of the calculus is a theorem. It follows that if $\sim A$ is not a theorem, then $\sim A$ is not valid; and, therefore, A is satisfiable; *i.e.*, it has an interpretation, or a model. But to say that A is consistent means nothing other than that $\sim A$ is not a theorem. Hence, from the completeness, it follows that if A is consistent, then A is satisfiable. Therefore, the semantic concepts of validity and satisfiability are seen to coincide with the syntactic concepts of derivability and consistency.

The Löwenheim-Skolem theorem. A finding closely related to the completeness theorem is the Löwenheim-Skolem theorem (1915, 1920), named after Leopold Löwenheim, a German schoolteacher, and Skolem, which says that if a sentence (or a formal system) has any model, it has a countable or enumerable model (*i.e.*, a model whose members can be matched with the positive integers). In the most direct method of proving this theorem, the logician is provided with very useful tools in model theory and in studies on relative consistency and independence in set theory.

In the predicate calculus there are certain reduction or normal-form theorems. One useful example is the prenex normal form: every sentence can be reduced to an equivalent sentence expressed in the prenex form—*i.e.*, in a form such that all the quantifiers appear at the beginning. This form is especially useful for displaying the central ideas of some of the proofs of the Löwenheim-Skolem theorem.

As an illustration, one may consider a simple schema in prenex form. "For every x , there is some y such that x bears the (arbitrary) relation M to y "; *i.e.*,

$$(\forall x)(\exists y)Mxy. \tag{3}$$

If (3) now has a model with a nonempty domain D , then, by a principle from set theory (the axiom of choice), there exists a function f of x , written $f(x)$, that singles out for each x a corresponding y . Hence, "For every x , x bears the relation M to $f(x)$ "; *i.e.*,

$$(\forall x)Mxf(x). \tag{4}$$

If a is now any object in D , then the countable subdomain ($a, f(a), f[f(a)], \dots$) already contains enough objects to satisfy (4) and therefore to satisfy (3). Hence, if (3) has any model, it has a countable model, which is in fact a submodel of the original.

An alternative proof, developed by Skolem in 1922 to avoid appealing to the principles of set theory, has turned out to be useful also for establishing the completeness of the calculus. Instead of using the function f as before, a can be arbitrarily denoted by 1. Since equation (3) is true, there must be some object y such that the number 1 bears the relation M to y , or symbolically $M1y$; and one of these y 's may be called 2. When this process is repeated indefinitely, one obtains

$$M12; M12 \cdot M23; M12 \cdot M23 \cdot M34; \dots \tag{5}$$

all of which are true in the given model. The argument is elementary, because in each instance one merely argues from "There exists some y such that n is M of y "—*i.e.*, ($\exists y)Mny$ —to "Let one such y be $n + 1$." Consequently, every member of (5) is true in some model. It is then possible to infer that all members of (5) are simultaneously true in some model—*i.e.*, that there is some way of assigning truth values to all its atomic parts so that all members of (5) will be true. Hence, it follows that (3) is true in some countable model.

The completeness theorem. Gödel's original proof of the completeness theorem is closely related to the second proof above. Consideration may again be given to all the sentences in (5) that contain no more quantifiers. If they are all satisfiable, then, as before, they are simultaneously satisfiable and (3) has a model. On the other hand, if (3) has no model, some of its terms—say $M12 \cdot \dots \cdot M89$ —are not satisfiable; *i.e.*, their negations

Consistency and decidability of PC

Proof independent of set theory

are tautologies (theorems of the propositional calculus). Thus, $\sim M12 \vee \dots \vee \sim M89$ is a tautology, and this remains true if 1, 2, . . . , 9 are replaced by variables, such as r, s, \dots, z ; hence, $\sim Mrs \vee \dots \vee \sim Myz$, being a tautology expressed in the predicate calculus as usually formulated, is a theorem in it. It is then easy to use the usual rules of the predicate calculus to derive also the statement, "There exists an x such that, for every y , x is not M of y "; i.e., $(\exists x)(\forall y)\sim Mxy$. In other words, the negation of (3) is a theorem of the predicate calculus. Hence, if (3) has no model, then its negation is a theorem of the predicate calculus. And, finally, if a sentence is valid (i.e., if its negation has no model), then it is itself a theorem of the predicate calculus.

The undecidability theorem and reduction classes. Given the completeness theorem, it follows that the task of deciding whether any sentence is a theorem of the predicate calculus is equivalent to that of deciding whether any sentence is valid or whether its negation is satisfiable.

Turing's method of proving that this class of problems is undecidable is particularly suggestive. Once the concept of mechanical procedure was crystallized, it was relatively easy to find absolutely unsolvable problems—e.g., the halting problem, which asks for each Turing machine the question of whether it will ever stop, beginning with a blank tape. In other words, each Turing machine operates in a predetermined manner according to what is given initially on the (input) tape; we consider now the special case of a blank tape and ask the special question whether the machine will eventually stop. This infinite class of questions (one for each machine) is known to be unsolvable.

Turing's method shows that each such question about a single Turing machine can be expressed by a single sentence of the predicate calculus so that the machine will stop if and only if that sentence is not satisfiable. Hence, if there were a decision procedure of validity (or satisfiability) for all sentences of the predicate calculus, then the halting problem would be solvable.

In more recent years (1962), Turing's formulation has been improved to the extent that all that is needed are sentences of the relatively simple form $(\forall x)(\exists y)(\forall z)Mxyz$, in which all the quantifiers are at the beginning; i.e., M contains no more quantifiers. Hence, given the unsolvability of the halting problem, it follows that, even for the simple class of sentences in the predicate calculus having the quantifiers $\forall \exists \forall$, the decision problem is unsolvable. Moreover, the method of proof also yields a procedure by which every sentence of the predicate calculus can be correlated with one in the simple form given above. Thus, the class of $\forall \exists \forall$ sentences forms a "reduction class." (There are also various other reduction classes.)

MODEL THEORY

Background and typical problems. In model theory one studies the interpretations (models) of theories formalized in the framework of formal logic, especially in that of the first-order predicate calculus with identity—i.e., in elementary logic. A first-order language is given by a collection S of symbols for relations, functions, and constants, which, in combination with the symbols of elementary logic, single out certain combinations of symbols as sentences. Thus, for example, in the case of the system N (see above *Example of a formal system*), the formation rules yield a language that is determined in accordance with a uniform procedure by the set (indicated by braces) of uninterpreted extralogical symbols:

$$L = \{S, +, \cdot, 0, 1\}.$$

A first-order theory is determined by a language and a set of selected sentences of the language—those sentences of the theory that are, in an arbitrary, generalized sense, the "true" ones (called the "distinguished elements" of the set). In the particular case of the system N , one theory T_a is built up on the basis of the language and the set of theorems of N , and another theory T_b is determined by the true sentences of N according to the natural interpretation or meaning of its language. In general, the language of N and any set of sentences of the language can be used to make up a theory.

Satisfaction of a theory by a structure: finite and infinite models. A realization of a language (for example, the one based on L) is a structure \mathfrak{A} identified by the six elements so arranged

$$\mathfrak{A} = \langle A, S_{\mathfrak{A}}, +_{\mathfrak{A}}, \cdot_{\mathfrak{A}}, 0_{\mathfrak{A}}, 1_{\mathfrak{A}} \rangle,$$

in which the second term is a function that assigns a member of the set A to each member of the set A , the next two terms are functions correlating each member of the Cartesian product $A \times A$ (i.e., from the set of ordered pairs $\langle a, b \rangle$ such that a and b belong to A) with a member of A , and the last two terms are members of A . The structure \mathfrak{A} satisfies, or is a model of, the theory T_a (or T_b) if all of the distinguished sentences of T_a (or T_b) are true in \mathfrak{A} (or satisfied by \mathfrak{A}). Thus, if \mathfrak{A} is the structure of the ordinary nonnegative integers $\langle \omega, S, +, \cdot, 0, 1 \rangle$, in which ω is the set of all such integers and $S, +, \cdot, 0$, and 1 the elements for their generation, then it is not only a realization of the language based on L but also a model of both T_a and T_b . Gödel's incompleteness theorem permits nonstandard models of T_a that contain more objects than ω but in which all the distinguished sentences of T_a (namely, the theorems of the system N) are true. Skolem's constructions (related to ultraproducts, see below) yield nonstandard models for both theory T_a and theory T_b .

The use of the relation of satisfaction, or being-a-model-of, between a structure and a theory (or a sentence) can be traced to the book *Wissenschaftslehre* (1837; *Theory of Science*) by Bernhard Bolzano, a Bohemian theologian and mathematician, and, in a more concrete context, to the introduction of models of non-Euclidean geometries about that time. In the mathematical treatment of logic, these concepts can be found in works of the late 19th-century German mathematician Ernst Schröder and in Löwenheim (in particular, in his paper of 1915). The basic tools and results achieved in model theory—such as the Löwenheim-Skolem theorem, the completeness theorem of elementary logic, and Skolem's construction of nonstandard models of arithmetic—were developed during the period from 1915 to 1933. A more general and abstract study of model theory began after 1950, in the work of Tarski and others.

One group of developments may be classified as refinements and extensions of the Löwenheim-Skolem theorem. These developments employ the concept of a "cardinal number," which—for a finite set—is simply the number at which one stops in counting its elements. For infinite sets, however, the elements must be matched from set to set instead of being counted, and the "sizes" of these sets must thus be designated by transfinite numbers. A rather direct generalization can be drawn that says that, if a theory has any infinite model, then, for any infinite cardinal number, it has a model of that cardinality. It follows that no theory with an infinite model can be categorical or such that any two models of the theory are isomorphic (i.e., matchable in one-to-one correspondence), because models of different cardinalities can obviously not be so matched. A natural question is whether a theory can be categorical in certain infinite cardinalities—i.e., whether there are cardinal numbers such that any two models of the theory of the same cardinality are isomorphic. According to a central discovery made in 1963 by the American mathematician Michael Morley, if a theory is categorical in any uncountable cardinality (i.e., any cardinality higher than the countable), then it is categorical in every uncountable cardinality. On the other hand, examples are known for all four combinations of countable and uncountable cardinalities: specifically, there are theories that are categorical (1) in every infinite cardinality, (2) in the countable cardinality but in no uncountable cardinality, (3) in every uncountable cardinality but not in the countable, and (4) in no infinite cardinality.

In another direction, there are "two-cardinal" problems that arise from the possibilities of changing, from one model to another, not only the cardinality of the domain of the first model but also the cardinality of some chosen property (such as being a prime number). Various answers to these questions have been found, including proofs of independence (based on the ordinary axioms employed in

Integers as a structure

Categorical theories in various cardinalities

Turing's undecidability proof

set theory) and proofs of conditional theorems made on the basis of certain familiar hypotheses of set theory.

Elementary logic. An area that is perhaps of more philosophical interest is that of the nature of elementary logic itself. On the one hand, the completeness discoveries seem to show in some sense that elementary logic is what the logician naturally wishes to have. On the other hand, he is still inclined to ask whether there might be some principle of uniqueness according to which elementary logic is the only solution that satisfies certain natural requirements on what a logic should be. The development of model theory has led to a more general outlook that enabled the Swedish logician Per Lindström to prove in 1969 a general theorem to the effect that, roughly speaking, within a broad class of possible logics, elementary logic is the only one that satisfies the requirements of axiomatizability and of the Löwenheim-Skolem theorem. Although Lindström's theorem does not settle satisfactorily whether or not elementary logic is the right logic, it does seem to suggest that mathematical findings can help the logician to clarify his concepts of logic and of logical truth.

Com-
pounded
models

A particularly useful tool for obtaining new models from the given models of a theory is the construction of a special combination called the "ultraproduct" of a family of structures (see below *Ultrafilters, ultraproducts, and ultrapowers*)—in particular, the ultrapower when the structures are all copies of the same structure (just as the product of a_1, \dots, a_n is the same as the power a^n , if $a_i = a$ for each i). The intuitive idea in this method is to establish that a sentence is true in the ultraproduct if and only if it is true in "almost all" of the given structures (*i.e.*, "almost everywhere"—an idea that was present in a different form in Skolem's construction of a nonstandard model of arithmetic in 1933). It follows that, if the given structures are models of a theory, then their ultraproduct is such a model also, because every sentence in the theory is true everywhere (which is a special case of "almost everywhere" in the technical sense employed). Ultraproducts have been applied, for example, to provide a foundation for what is known as "nonstandard analysis" that yields an unambiguous interpretation of the classical concept of infinitesimals—the division into units as small as one pleases. They have also been applied by two mathematicians, James Ax and Simon B. Kochen, to problems in the field of algebra (on p -adic fields).

Nonelementary logic and future developments. There are also studies, such as second-order logic and infinitary logics, that develop the model theory of nonelementary logic. Second-order logic contains, in addition to variables that range over individual objects, a second kind of variable ranging over sets of objects so that the model of a second-order sentence or theory also involves, beyond the basic domain, a larger set (called its "power set") that encompasses all the subsets of the domain. Infinitary logics may include functions or relations with infinitely many arguments, infinitely long conjunctions and disjunctions, or infinite strings of quantifiers. From studies on infinitary logics, William Hanf, an American logician, was able to define certain cardinals, some of which have been studied in connection with the large cardinals in set theory. In yet another direction, logicians are developing model theories for modal logics—those dealing with such modalities as necessity and possibility—and for the intuitionistic logic.

There is a large gap between the general theory of models and the construction of interesting particular models such as those employed in the proofs of the independence (and consistency) of special axioms and hypotheses in set theory. It is natural to look for further developments of model theory that will yield more systematic methods for constructing models of axioms with interesting particular properties, especially in deciding whether certain given sentences are derivable from the axioms. Relative to the present state of knowledge, such goals appear fairly remote. The gap is not unlike that between the abstract theory of computers and the basic properties of actual computers.

Characterizations of the first-order logic. There has been outlined above a proof of the completeness of elementary logic without including sentences asserting identity. The proof can be extended, however, to the full

elementary logic in a fairly direct manner. Thus, if F is a sentence containing equality, a sentence G can be adjoined to it that embodies the special properties of identity relevant to the sentence F . The conjunction of F and G can then be treated as a sentence not containing equality (*i.e.*, "=" can be treated as an arbitrary relation symbol). Hence, the conjunction has a model in the sense of logic-without-identity if and only if F has a model in the sense of logic-with-identity; and the completeness of elementary logic (with identity) can thus be inferred.

A concept more general than validity is that of the relation of logical entailment or implication between a possibly infinite set X of sentences and a single sentence p that holds if and only if p is true in every model of X . In particular, p is valid if the empty set, defined as having no members, logically entails p —for this is just another way of saying that p is true in every model. This suggests a stronger requirement on a formal system of logic—namely, that p be derivable from X by the system whenever X logically entails p . The usual systems of logic satisfy this requirement because, besides the completeness theorem, there is also a compactness theorem:

A theory X has a model if every finite subset of X has a model.

Roughly speaking, this theorem enables the logician to reduce an infinite set X to a finite subset X_1 in each individual case, and the case of entailment when X_1 is finite is taken care of by the completeness of the system.

These findings show that the ordinary systems of elementary logic comprise the correct formulation, provided that the actual choice of the truth functions (say negation and disjunction), of the quantifiers, and of equality as the "logical constants" is assumed to be the correct one. There remains the question, however, of justifying the particular choice of logical constants. One might ask, for example, whether "For most x " or "For finitely many x " should not also be counted as logical constants. Lindström has formulated a general concept of logic and shown that logics that apparently extend the first-order logic all end up being the same as that logic, provided that they satisfy the Löwenheim-Skolem theorem and either have the compactness property or are formally axiomatizable. There remains the question, however, of whether or why these requirements (especially that of the Löwenheim-Skolem theorem) are intrinsic to the nature of logic.

Generalizations and extensions of the Löwenheim-Skolem theorem. A generalized theorem can be proved using basically the same ideas as those employed in the more special case discussed above.

If a theory has any infinite model, then, for any infinite cardinality α , that theory has a model of cardinality α . More explicitly, this theorem contains two parts: (1) If a theory has a model of infinite cardinality β , then, for each infinite cardinal α that is greater than β , the theory has a model of cardinality α . (2) If a theory has a model of infinite cardinality β , then, for each infinite cardinal α less than β , the theory has a model of cardinality α .

It follows immediately that any theory having an infinite model has two nonisomorphic models and is, therefore, not categorical. This applies, in particular, to the aforementioned theories T_a and T_b of arithmetic (based on the language of \mathbb{N}), the natural models of which are countable, as well as to theories dealing with real numbers and arbitrary sets, the natural models of which are uncountable; both kinds of theory have both countable and uncountable models. There is much philosophical discussion about this phenomenon.

The possibility is not excluded that a theory may be categorical in some infinite cardinality. The theory T_b , for example, of dense linear ordering (such as that of the rational numbers) is categorical in the countable cardinality. One application of the Löwenheim-Skolem theorem is: If a theory has no finite models and is categorical in some infinite cardinality α , then the theory is complete; *i.e.*, for every closed sentence in the language of the theory, either that sentence or its negation belongs to the theory. An immediate consequence of this application of the theorem is that the theory of dense linear ordering is complete.

Extension
to include
identity

Morley's theorem

A theorem that is generally regarded as one of the most difficult to prove in model theory is the theorem by Michael Morley, as follows:

A theory that is categorical in one uncountable cardinality is categorical in every uncountable cardinality.

Two-cardinal theorems deal with languages having some distinguished predicate U . A theory is said to admit the pair $\langle \alpha, \beta \rangle$ of cardinals if it has a model (with its domain) of cardinality α wherein the value of U is a set of cardinality β . The central two-cardinal theorem says:

If a theory admits the pair $\langle \alpha, \beta \rangle$ of infinite cardinals with β less than α , then for each regular cardinal γ the theory admits $\langle \gamma^+, \gamma \rangle$, in which γ^+ is the next larger cardinal after γ .

The most interesting case is when γ is the least infinite cardinal, \aleph_0 . (The general theorem can be established only when the "generalized continuum hypothesis" is assumed, according to which the next highest cardinality for an infinite set is that of its power set.)

Ultrafilters, ultraproducts, and ultrapowers. An ultrafilter on a nonempty set I is defined as a set D of subsets of I such that

- (1) the empty set does not belong to D ,
- (2) if A, B are in D , so is their intersection, $A \cap B$, the set of elements common to both,
- (3) if A is a subset of B , and A is in D , then B is in D , and
- (4) for every subset A of I , either A is in D or I minus A is in D .

Roughly stated, each ultrafilter of a set I conveys a notion of large subsets of I so that any property applying to a member of D applies to I "almost everywhere."

The set $\{\mathfrak{A}_i\}$, where $\mathfrak{A}_i = \langle A_i, R_i \rangle$ and the i are members of the set I , is taken to be a family of structures indexed by I , and D to be an ultrafilter on I . Consider now the Cartesian product B of $\{A_i\}$ (for example, if I is $\{0, 1, 2, \dots\}$, then B is the set of all sequences f such that $f(i)$ belong to A_i). The members of B are divided into equivalence classes with the help of $D: f \equiv g$ if and only if $\{i \mid f(i) = g(i)\} \in D$ —in other words, the set of indices i such that $f(i) = g(i)$ belong to D [or $f(i)$ and $g(i)$ are equal "almost everywhere"]. Let W be the set of these equivalence classes—*i.e.*, the set of all f^* such that f^* is the set of all members g of B with $g \equiv f$. Similarly, a relation S is introduced such that Sfg if and only if R_i holds between $f(i)$ and $g(i)$ for "almost all" i ; *i.e.*,

$$\{i \mid R_i [f(i), g(i)]\} \in D.$$

In this way, we arrive at a new structure $U = \langle W, S \rangle$, which is called the ultraproduct of the original family $\{\mathfrak{A}_i\}$ over D . In the special case when all the \mathfrak{A}_i are the same, the resulting structure U is called the ultrapower of the original family over D .

The central theorems are the following:

1. If \mathfrak{A}_i ($i \in I$) are realizations of the same language, then a sentence p is true in the ultraproduct U if and only if the set of i such that p is true in \mathfrak{A}_i belongs to D . In particular, if each \mathfrak{A}_i is a model of a theory, then U is also a model of the theory.
2. Two realizations of the same language are said to be elementarily equivalent if they have the same set of true sentences. A necessary and sufficient condition for two realizations to be elementarily equivalent is that they admit ultrapowers that are isomorphic.

One application of these theorems is in the introduction of nonstandard analysis, which was originally instituted by other considerations. By using a suitable ultrapower of the structure of the field \mathfrak{R} of real numbers, a real closed field that is elementarily equivalent to \mathfrak{R} is obtained that is non-Archimedean—*i.e.*, which permits numbers a and b such that no n can make na greater than b . This development supplies an unexpected exact foundation for the classical differential calculus using infinitesimals, which has considerable historical, pedagogical, and philosophical interest.

A widely known application to the area of algebra is that which deals with certain fields of rational numbers Q_p , called the p -adic completion of the rational numbers. The conjecture has been made that every form of degree

d (in the same sense as degrees of ordinary polynomials) over Q_p , in which the number of variables exceeds d^2 , has a nontrivial zero in Q_p . Using ultraproducts, it has been shown that the conjecture is true for arbitrary d with the possible exception of a finite set of primes p (depending on d). Subsequently, it was found that the original conjecture is not true when extended to full generality.

Other useful tools in model theory include the pigeon-hole principles, of which the basic principle is that, if a set of large cardinality is partitioned into a small number of classes, some one class will have large cardinality. Those elements of the set that lie in the same class cannot be distinguished by the property defining that class.

A related concept is that of "indiscernibles," which also has rather extensive applications in set theory. An ordered subset of the domain of a model \mathfrak{A} of a theory is a homogeneous set, or a set of indiscernibles for \mathfrak{A} , if \mathfrak{A} cannot distinguish the members of the subset from one another. More exactly, given any $x_1 < \dots < x_n, y_1 < \dots < y_n$ in the subset, then for any sentence $F(a_1, \dots, a_n)$ of the language of the theory, that sentence (with argument x) is satisfied by (symbolized \models) the structure—*i.e.*,

$$\mathfrak{A} \models F(x_1, \dots, x_n)$$

—if and only if that sentence (with argument y) is also satisfied by it—*i.e.*,

$$\mathfrak{A} \models F(y_1, \dots, y_n).$$

There is also a first theorem on this notion that says that, given a theory with an infinite model and a linearly ordered set X , there is then a model \mathfrak{A} of the theory such that X is a set of indiscernibles for \mathfrak{A} . (H.Wa./M.L.Sc.)

Applied logic

The formalism and theoretical results of pure logic can be clothed with meanings derived from a variety of sources within philosophy as well as from other sciences. This formal machinery also can be used to guide the design of computers and computer programs.

The applications of logic cover a vast range, relating to reasoning in the sciences and in philosophy, as well as in everyday discourse. They include (1) the various sorts of reasoning affecting the conduct of ordinary discourse as well as the theory of the logical relations that exist within special realms of discourse—between two commands, for example, or between one question and another, (2) special forms of logic designed for scientific applications, such as temporal logic (of what "was" or "will be" the case) or mereology (the logic of parts and wholes), and (3) special forms for concepts bearing upon philosophical issues, such as logics that deal with statements of the form "I know that . . .," "I believe that . . .," "It is permitted to . . .," "It is obligatory to . . .," or "It is prohibited to . . ."

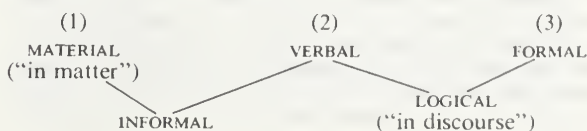
Range of applications

THE CRITIQUE OF FORMS OF REASONING

Correct and defective argument forms. In logic an argument consists of a set of statements, the premises, whose truth supposedly supports the truth of a single statement called the conclusion of the argument. An argument is deductively valid when the truth of the premises guarantees the truth of the conclusion; *i.e.*, the conclusion must be true, because of the form of the argument, whenever the premises are true. Some arguments that fail to be deductively valid are acceptable on grounds other than formal logic, and their conclusions are supported with less than logical necessity. In other potentially persuasive arguments, the premises give no rational grounds for accepting the conclusion. These defective forms of argument are called fallacies.

An argument may be fallacious in three ways: in its material content, through a misstatement of the facts; in its wording, through an incorrect use of terms; or in its

Applications to non-standard analysis and algebra



structure (or form), through the use of an improper process of inference. As shown in the diagram, fallacies are correspondingly classified as (1) material, (2) verbal, and (3) formal. Groups 2 and 3 are called logical fallacies, or fallacies "in discourse," in contrast to the substantive, or material, fallacies of group 1, called fallacies "in matter"; and groups 1 and 2, in contrast to group 3, are called informal fallacies.

Kinds of fallacies. *Material fallacies.* The material fallacies are also known as fallacies of presumption, because the premises "presume" too much—they either covertly assume the conclusion or avoid the issue in view.

The classification that is still widely used is that of Aristotle's *Sophistic Refutations*: (1) The fallacy of accident is committed by an argument that applies a general rule to a particular case in which some special circumstance ("accident") makes the rule inapplicable. The truth that "men are capable of seeing" is no basis for the conclusion that "blind men are capable of seeing." This is a special case of the fallacy of *secundum quid* (more fully: *a dicto simpliciter ad dictum secundum quid*, which means "from a saying [taken too] simply to a saying according to what [it really is]"—i.e., according to its truth as holding only under special provisos). This fallacy is committed when a general proposition is used as the premise for an argument without attention to the (tacit) restrictions and qualifications that govern it and invalidate its application in the manner at issue. (2) The converse fallacy of accident argues improperly from a special case to a general rule. Thus, the fact that a certain drug is beneficial to some sick persons does not imply that it is beneficial to all people. (3) The fallacy of irrelevant conclusion is committed when the conclusion changes the point that is at issue in the premises. Special cases of irrelevant conclusion are presented by the so-called fallacies of relevance. These include (a) the argument *ad hominem* (speaking "against the man" rather than to the issue), in which the premises may only make a personal attack on a person who holds some thesis, instead of offering grounds showing why what he says is false, (b) the argument *ad populum* (an appeal "to the people"), which, instead of offering logical reasons, appeals to such popular attitudes as the dislike of injustice, (c) the argument *ad misericordiam* (an appeal "to pity"), as when a trial lawyer, rather than arguing for his client's innocence, tries to move the jury to sympathy for him, (d) the argument *ad verecundiam* (an appeal "to awe"), which seeks to secure acceptance of the conclusion on the grounds of its endorsement by persons whose views are held in general respect, (e) the argument *ad ignorantiam* (an appeal "to ignorance"), which argues that something (e.g., extrasensory perception) is so since no one has shown that it is not so, and (f) the argument *ad baculum* (an appeal "to force"), which rests on a threatened or implied use of force to induce acceptance of its conclusion. (4) The fallacy of circular argument, known as *petitio principii* ("begging the question"), occurs when the premises presume, openly or covertly, the very conclusion that is to be demonstrated (example: "Gregory always votes wisely." "But how do you know?" "Because he always votes Libertarian."). A special form of this fallacy, called a vicious circle, or *circulus in probando* ("arguing in a circle"), occurs in a course of reasoning typified by the complex argument in which a premise p_1 is used to prove p_2 ; p_2 is used to prove p_3 ; and so on, until p_{n-1} is used to prove p_n ; then p_n is subsequently used in a proof of p_1 , and the whole series p_1, p_2, \dots, p_n is taken as established (example: "McKinley College's baseball team is the best in the association [$p_n = p_3$]; they are the best because of their strong batting potential [p_2]; they have this potential because of the ability of Jones, Crawford, and Randolph at the bat [p_1]." "But how do you know that Jones, Crawford, and Randolph are such good batters?" "Well, after all, these men are the backbone of the best team in the association [p_3 again]."). Strictly speaking, *petitio principii* is not a fallacy of reasoning but an ineptitude in argumentation: thus the argument from p as a premise to p as conclusion is not deductively invalid but lacks any power of conviction, since no one who questioned the conclusion could concede the premise. (5)

Accident,
converse
accident,
irrelevant
conclusion

Circular
argument

The fallacy of false cause (*non causa pro causa*) mislocates the cause of one phenomenon in another that is only seemingly related. The most common version of this fallacy, called *post hoc ergo propter hoc* ("after which hence by which"), mistakes temporal sequence for causal connection—as when a misfortune is attributed to a "malign event," like the dropping of a mirror. Another version of this fallacy arises in using *reductio ad absurdum* reasoning: concluding that a statement is false if its addition to a set of premises leads to a contradiction. This mode of reasoning can be correct—e.g., concluding that two lines do not intersect if the assumption that they do intersect leads to a contradiction. What is required to avoid the fallacy is to verify independently that each of the original premises is true. Thus, one might fallaciously infer that Williams, a philosopher, does not watch television, because adding

A: Williams, a philosopher, watches television.

to the premises

P_1 : No philosopher engages in intellectually trivial activities.

P_2 : Watching television is an intellectually trivial activity.

leads to a contradiction. Yet it might be that either P_1 or P_2 or both are false. It might even be the case that Williams is not a philosopher. Indeed, one might even take A as evidence for the falsity of either P_1 or P_2 or as evidence that Williams is not really a philosopher. (6) The fallacy of many questions (*plurimum interrogationum*) consists in demanding or giving a single answer to a question when this answer could either be divided (example: "Do you like the twins?" "Neither yes nor no; but Ann yes and Mary no.") or refused altogether, because a mistaken presupposition is involved (example: "Have you stopped beating your wife?"). (7) The fallacy of *non sequitur* ("it does not follow") occurs when there is not even a deceptively plausible appearance of valid reasoning, because there is an obvious lack of connection between the given premises and the conclusion drawn from them. Some authors, however, identify *non sequitur* with the fallacy of the consequent (see below *Formal fallacies*).

Verbal fallacies. These fallacies, called fallacies of ambiguity, arise when the conclusion is achieved through an improper use of words. The principal instances are as follows: (1) Equivocation occurs when a word or phrase is used in one sense in one premise and in another sense in some other needed premise or in the conclusion (example: "The loss made Jones mad [= angry]; mad [= insane] people should be institutionalized; so Jones should be institutionalized."). The figure-of-speech fallacy is the special case arising from confusion between the ordinary sense of a word and its metaphorical, figurative, or technical employment (example: "For the past week Joan has been living on the heights of ecstasy." "And what is her address there?"). (2) Amphiboly occurs when the grammar of a statement is such that several distinct meanings can obtain (example: "The governor says, 'Save soap and waste paper.' So soap is more valuable than paper"). (3) Accent is a counterpart of amphiboly arising when a statement can bear distinct meanings depending on which word is stressed (example: "Men are considered equal." "Men are considered equal."). (4) Composition occurs when the premise that the parts of a whole are of a certain nature is improperly used to infer that the whole itself must also be of this nature (example: a story made up of good paragraphs is thus said to be a good story). (5) Division—the reverse of composition—occurs when the premise that a collective whole has a certain nature is improperly used to infer that a part of this whole must also be of this nature (example: in a speech that is long-winded it is presumed that every sentence is long). But this fallacy and its predecessor can be viewed as versions of equivocation, in which the distributive use of a term—i.e., its application to the elements of an aggregate (example: "the crowd," viewed as individuals)—is confused with its collective use ("the crowd," as a unitary whole)—compare "The crowd were filing through the turnstile" with "The crowd was compressed into the space of a city block."

Other
material
fallacies

Equivoca-
tion,
amphiboly,
and other
verbal
fallacies

Formal fallacies. Formal fallacies are deductively invalid arguments that typically commit an easily recognizable logical error. A classic case is Aristotle's fallacy of the consequent, relating to reasoning from premises of the form "If p_1 , then p_2 ." The fallacy has two forms: (1) denial of the antecedent, in which one mistakenly argues from the premises "If p_1 , then p_2 " and "not- p_1 " (symbolized $\sim p_1$) to the conclusion "not- p_2 ." (example: "If George is a man of good faith, he can be entrusted with this office; but George is not a man of good faith; therefore, George cannot be entrusted with this office"), and (2) affirmation of the consequent, in which one mistakenly argues from the premises "If p_1 , then p_2 " and " p_2 " to the conclusion " p_1 ." (example: "If Amos was a prophet, then he had a social conscience; he had a social conscience; hence, Amos was a prophet"). Most of the traditionally considered formal fallacies, however, relate to the syllogism. One example may be cited, that of the fallacy of illicit major (or minor) premise, which violates the rules for "distribution." (A term is said to be distributed when reference is made to all members of the class. For example, in "Some crows are not friendly," reference is made to all friendly things but not to all crows.) The fallacy arises when a major (or minor) term that is undistributed in the premise is distributed in the conclusion (example: "All tubers are high-starch foods [undistributed]; no squashes are tubers; therefore, no squashes are high-starch foods [distributed]").

EPISTEMIC LOGIC

Epistemic logic deals with the logical issues arising within the gamut of such epistemological concepts as knowledge, belief, assertion, doubt, question-and-answer, or the like. Instead of dealing with the essentially factual issues of alethic logic (Greek: *alētheia*, "truth")—*i.e.*, with what is actually or must necessarily or can possibly be the case—it relates to what people know or believe or maintain or doubt to be the case.

The logic of belief. From the logical standpoint, a belief is generally analyzed as a relationship obtaining between the person who accepts some thesis on the one hand and the thesis that he accepts on the other. Correspondingly, given a person x , it is convenient to consider the set B_x of x 's beliefs and represent the statement " x believes that p " as $p \in B_x$. (The symbol \in represents membership in a set, \notin its denial.)

Theory of belief. To articulate a viable logic of belief, it is, at the very least, essential to postulate certain minimal conditions of rationality regarding the parties whose beliefs are at issue:

1. Consistency: "If x believes that p , then x does not believe that not- p "; *i.e.*,

$$\text{If } p \in B_x, \text{ then } \sim p \notin B_x.$$

"If not- p , then x does not believe that p "; *i.e.*,

$$\text{If } \sim p, \text{ then } p \notin B_x.$$

Example: If "Jesus was a Zealot" (p) is among (\in) the beliefs of Ralph (B_{Ralph}), then "Jesus was not a Zealot" ($\sim p$) is not among (\notin) Ralph's beliefs. It is an accepted thesis (\vdash) that "Jesus was not a Zealot." Hence, "Jesus was a Zealot" is not among Ralph's beliefs. (The symbol " \vdash " is used to indicate that the sentence to its right is a valid deductive consequence of the sentence[s] on the left. In cases where it appears as an isolated prefix, it signifies "theoremhood"—*i.e.*, a deductive consequence from no premises.)

2. Conjunctive composition and division: "If x believes that p_1 , and x believes that p_2 , etc., to x believes that p_n , then x believes that p_1 and p_2 , etc., and p_n "; *i.e.*,

$$\text{If } (p_1 \in B_x, p_2 \in B_x, \dots, p_n \in B_x), \\ \text{then } (p_1 \cdot p_2 \cdot \dots \cdot p_n) \in B_x,$$

and conversely. Example: If "cats are affectionate" (p_1), "cats are clean" (p_2), etc., to "cats are furry" (p_n) are among (\in) Bob's beliefs (B_{Bob}), then "cats are affectionate and clean, etc., and furry" ($p_1 \cdot p_2 \cdot \dots \cdot p_n$) is also a belief of Bob's.

3. Minimal inferential capacity: "If x believes that p ,

and q is an obvious consequence of p , then x believes that q "; *i.e.*,

$$\text{If } p \in B_x \text{ and } p \vDash q, \text{ then } q \in B_x.$$

Example: "If x believes that his cat is on the mat, and his cat's being on the mat has an obvious consequence that something is on the mat, then x believes that something is on the mat."

Here item 3 is a form of the entailment principle, but with \vDash representing entailment of the simplest sort, designating obvious consequence—say, deducibility by fewer than two (or n) inferential steps, employing only those primitive rules of inference that have been classified as obvious. (In arguments about beliefs, however, all repetitions of the application of this version of the entailment principle must be avoided.) These principles endow the theory with such rules as

1. "If x believes that not- p , then x does not believe that p "; *i.e.*,

$$\text{If } \sim p \in B_x, \text{ then } p \notin B_x.$$

2. "If x believes that p , and x believes that q , then x believes that both p and q taken together"; *i.e.*,

$$\text{If } p \in B_x \text{ and } q \in B_x, \text{ then } p \cdot q \in B_x.$$

3. "If x believes that p , then x believes that either p or q "; *i.e.*,

$$\text{If } p \in B_x, \text{ then } p \vee q \in B_x,$$

given " $p \vdash p \vee q$ " as an "obvious" rule of inference (where \vee means "or").

One key question of the logical theory of belief relates to the area of iterative beliefs (example: "Andrews believes that I believe that he believes me to be untrustworthy"). Clearly, one would not want to have such theses as:

Iterative beliefs

1. "If y believes that x believes that p , then x believes that p "; *i.e.*,

$$\text{If } (p \in B_x) \in B_y, \text{ then } p \in B_x \quad (y \neq x)$$

2. "If y believes that x believes that p , then y believes that p "; *i.e.*,

$$\text{If } (p \in B_x) \in B_y, \text{ then } p \in B_y \quad (y \neq x)$$

But when the iteration is subject-uniform rather than subject-diverse, it might be advantageous to postulate certain special theses, such as

$$\text{If } p \in B_x, \text{ then } (p \in B_x) \in B_x,$$

which in effect limits the beliefs at issue to conscious beliefs. The plausibility of this thesis also implicates its converse—namely, whether there are circumstances under which someone's believing that he believes something would necessarily vouch for his believing of it (that is, whether it is legitimate to argue that "if x believes that he believes that p , then he believes that p "); *i.e.*,

$$\text{If } (p \in B_x) \in B_x, \text{ then } p \in B_x.$$

According to this thesis, the belief set B_x is to have the feature of second-order—as opposed to direct—applicability. From $q \in B_x$, it is not, in general, permissible to infer q , but one is entitled to do so when q takes the special form $p \in B_x$ —*i.e.*, when the belief at issue is one about the subject's own beliefs.

The theory is predicated on the view that belief is subject to logical compulsion but that the range of this compulsion is limited since people are not logically omniscient. Belief here is like sight: man has a limited range of logical vision; he can see clearly in the immediate logical neighbourhood of his beliefs but only dimly afar.

The logic of knowing. The propositional sense of knowing (*i.e.*, knowing that something or other is the case), rather than the operational sense of knowing (*i.e.*, knowing how something or other is done), is generally taken as the starting point for a logical theory of knowing. Accordingly, the logician may begin with a person x and consider a set of propositions K_x to represent his "body of knowledge." The aim of the theory then is to clarify and to characterize the relationship " x knows that p " or " p is among the items known to x ," which is here represented as $p \in K_x$.

Conditions of rationality

“Knowing as” and “being” true

There can be false knowledge only in the sense that “he thought he knew that p , but he was mistaken.” When the falsity of purported knowledge becomes manifest, the claim to knowledge must be withdrawn. “I know that p , but it may be false that p ” is a contradiction in terms. When something is asserted or admitted as known, it follows that this must be claimed to be true. But what sort of inferential step is at issue in the thesis that “ x knows p ” leads to “ p is true”? Is the link deductive, inductive, presuppositional, or somehow “pragmatic”? Each view has its supporters: on the deductive approach, $p \in K_x$ logically implies (deductively entails) p ; on the inductive approach, $p \in K_x$ renders p extremely probable, though not necessarily certain; on the presuppositional approach, $p \in K_x$ is improper (nonsensical) whenever p is not true; and on the pragmatic approach, the assertion of $p \in K_x$ carries with it a rational commitment to the assertion of p (in a manner, however, that does not amount to deductive entailment). From the standpoint of a logic of knowing, the most usual practice is to assume the deductive approach and to lay it down as a rule that if $p \in K_x$, then p is true. This approach construes knowledge in a very strong sense.

According to a common formula, knowledge is “true, justified belief.” This formulation, however, seems defective. Let the expression $J_x p$ be defined as meaning “ x has justification for accepting p ”; then

$$p \in K_x = p \cdot J_x p \cdot p \in B_x$$

For example, the proposition “Jane knows that (K_{Jane}) the gown is priceless (p)” means ($=$) “The gown is priceless, and Jane has justification for accepting that it is priceless ($J_{Jane} p$) and Jane believes that it is priceless ($p \in B_{Jane}$).” One cannot but assume that the conceptual nature of J is such as to underwrite the rule: “If x is justified in accepting p , then he is justified in accepting ‘Either p or q ’”; i.e.,

$$\text{If } J_x p, \text{ then } J_x(p \vee q), \quad (J)$$

in which q can be any other proposition whatsoever. The components p , q , and x may be such that all of the following obtain:

1. not- p
2. q
3. x believes that p ; i.e., $p \in B_x$
4. x does not believe that q ; i.e., $q \notin B_x$ and, indeed, x believes that not- q ; i.e., $\sim q \in B_x$
5. x is justified in accepting q ; i.e., $J_x q$
6. x believes that either p or q ; i.e., $p \vee q \in B_x$

Clearly, on any reasonable interpretation of B and J , this combination of six premises is possible. But the following consequences would then obtain:

7. $p \vee q$ by item 2 above
8. $J_x(p \vee q)$ by item 5 above and by J
9. $(p \vee q) \in K_x$ by items 6, 7, and 8

The conclusion (9) is wrong, however; x cannot properly be said to know that either p or q when $p \vee q$ is true solely because of the truth of q (which x rejects), but $p \vee q$ is believed by x solely because he accepts p (which is false). This example shows that the proposed definition of knowledge as “true, justified belief” cannot be made to work. The best plan, therefore, seems to be to treat the logic of knowing directly, rather than through the mediation of acceptance (belief) and justification.

Since Aristotle’s day, stress has been placed on the distinction between actual, overt knowledge that requires an explicit, consciously occurring awareness of what is known and potential, tacit knowledge that requires only implicit dispositional awareness. Unless $p \in K_x$ is construed in the tacit sense, the following principles will not hold:

$$\begin{aligned} &\text{If } p \in K_x \text{ and } p \vdash q, \text{ then } q \in K_x. \\ &\text{If } p \in K_x \text{ and } q \in K_x, \text{ then } (p \cdot q) \in K_x. \end{aligned}$$

These two rules, if accepted, however, suffice to guarantee the principle

$$\text{If } p_1, p_2, \dots, p_n \vdash q, \text{ then } p_1 \in K_x, p_2 \in K_x, \dots, p_n \in K_x \vdash q \in K_x.$$

Similar considerations regarding the potential construction of knowledge govern the answer to the question of whether, when something is known, this fact itself is

Overt versus tacit knowledge

known: if $p \in K_x$, then $(p \in K_x) \in K_x$. This principle is eminently plausible, provided that the membership of K_x is construed in the implicit (tacit) rather than in the explicit (overt) sense.

The logic of questions. Whether a given grouping of words is functioning as a question may hinge upon intonation, accentuation, or even context, rather than upon overt form: at bottom, questions represent a functional rather than a purely grammatical category. The very concept of a question is correlative with that of an answer, and every question correspondingly delimits a range of possible answers. One way of classifying questions is in terms of the surface characteristics of this range. On this basis, the logician can distinguish (among others):

- (1) yes/no questions (example: “Is today Tuesday?”),
- (2) item-specification questions (example: “What is an instance of a prime number?”),
- (3) instruction-seeking questions (example: “How does one bake an apple pie?”), and so on.

From the logical standpoint, however, a more comprehensive policy and one leading to greater precision is to treat every answer as given in a complete proposition (“Today is not Tuesday,” “Three is an example of a prime number,” and so on). From this standpoint, questions can be classed in terms of the nature of the answers. There would then be factual questions (example: “What day is today?”) and normative questions (example: “What ought to be done in these circumstances?”).

The advantage of the propositional approach to answers is that it captures the intrinsically close relationship between question and answer. The possible answers to (1) “What is the population of A-ville?” and (2) “What is the population of B-burgh?” are seemingly the same—namely, numbers of the series 0, 1, 2, But once complete propositions are taken to be at issue, then an answer to 1, such as “The population of A-ville is 5,238,” no longer counts as an answer to 2, since the latter must mention B-burgh. This approach has the disadvantage, on the other hand, of obscuring similarities in similar questions. One can now no longer say of two brothers that the questions “Who is Tom’s father?” and “Who is John’s father?” have the same answer.

With every question Q can be correlated the set of propositions $A(Q)$ of possible answers to Q . Thus, “What day of the week is today?” has seven conceivable answers, of the form “The day of the week today is Monday,” and the like. A possible answer to a question must be a possibly true statement. Accordingly, the question “What is an example of a prime number?” does not have “The Washington Monument is an example of a prime number” among its possible answers.

A question can be said to be true if it has a true answer—i.e., if $(\exists p) [p \cdot p \in A(Q)]$, which (taking the existential quantifier \exists to mean “there exists . . .”) can be read “There exists a proposition p such that p is true and p is among the answers of Q .” Otherwise it is false—i.e., all its answers are false. If he never came at all, the question “On what day of the week did he come?” is a false question in the sense that it lacks any true answer.

A true question can be called contingent if it admits of possible answers that are false, as in “Where did Jones put his pen?” In logic and mathematics there are, presumably, no contingent questions.

Questions can have presuppositions, as in “Why does Smith dislike Jones?” Any possible answer here must take the form “Smith dislikes Jones because . . .” and so commits one to the claim that “Smith dislikes Jones.” Every such question with a false presupposition must be a false question: all its possible answers (if any) are false.

Besides falsity, questions can exhibit an even more drastic sort of “impropriety.” They can be illegitimate in that they have no possible answers whatsoever (example: “What is an example of an even prime number different from two?”). The logic of questions is correspondingly three-valued: a question can be true (i.e., have a true answer), illegitimate (i.e., have no possible answer at all), or false (i.e., have possible answers but no true ones).

One question, Q_1 , will entail another, Q_2 , if every possible answer to the first deductively yields a possible answer to

Propositional approach

Truth-values in question theory

the second, and every true answer to the first deductively yields a true answer to the second. In this sense the question "What are the dimensions of that box?" entails the question "What is the height of that box?"

PRACTICAL LOGIC

The theory of reasoning with concepts of practice—of analyzing the logical relations obtaining among statements about actions and their accompaniments in choosing, planning, commanding, permitting, and so on—constitutes the domain of practical logic.

The logic of preference. The logic of preference—also called the logic of choice, or proairetic logic (Greek *proairesis*, "a choosing")—seeks to systematize the formal rules that govern the conception "x is preferred to y." A diversity of things can be at issue here: (1) Is x preferred to y by some individual (or group), or is x preferable to y in terms of some impersonal criterion? (2) Is on-balance preferability at issue or preferability in point of some particular factor (such as economy or safety or durability)? The resolution of these questions, though vital for interpretation, does not affect the formal structure of the preference relationships.

Symbolization and approach taken in proairetic logic. The fundamental tools of the logic of preference are as follows: (1) (strong) preference: x is preferable to y, symbolically $x \gg y$, (2) indifference: x and y are indifferent, $x \equiv y$, defined as "neither $x \gg y$ nor $y \gg x$," and (3) weak preference: x is no less preferred than y, $x \geq y$, defined as "either $x \gg y$ or $x \equiv y$." Since preference constitutes a relationship, its three types can be classed in terms of certain distinctions commonly drawn in the logic of relations: that of reflexivity (whether holding of itself: "John supports himself"), that of symmetry (whether holding when its terms are interchanged: "Peter is the cousin of Paul"; "Paul is the cousin of Peter"), and that of transitivity (whether transferable: $a \gg b$ and $b \gg c$; therefore $a \gg c$). Once it is established that the (strong) preference relation (\gg) is an ordering (i.e., is irreflexive, asymmetric, and transitive), it then follows that weak preference (\geq) is reflexive, nonsymmetric, and transitive and that indifference (\equiv) is an equivalence relation (i.e., reflexive, symmetric, and transitive).

One common approach to establishing a preference relation is to begin with a "measure of merit" to evaluate the relative desirability of the items x, y, z, . . . , that are at issue. Thus for any item x, a real-number quantity is obtained, symbolized # (x). (Such a measure is called a utility measure, the units are called utiles, and the comparisons or computations involved constitute a preference calculus.) In terms of such a measure, a preference ordering is readily introduced by the definitions that (1) $x \gg y$ is to be construed as $\#(x) > \#(y)$, (2) $x \geq y$ as $\#(x) \geq \#(y)$, and (3) $x \equiv y$ as $\#(x) = \#(y)$, in which \geq means "is greater than or equal to." Given these definitions, the relationships enumerated above must all obtain. Thus, the step from a utility measure to a preference ordering is simple.

Construction of a logic of preference. In constructing a logic of preference, it is assumed that the items at issue are propositions p, q, r, . . . and that the logician is to introduce a preferential ordering among them, with $p \gg q$ to mean "p's being the case is preferred to q's being the case." The problem is to systematize the logical relationships among such statements in order to permit a determination of whether, for example, it is acceptable to argue that "if either p is preferable to q or p is preferable to r, then p is preferable to either q or r," symbolized

$$(p \gg q \vee p \gg r) \supset [p \gg (q \vee r)]$$

(in which \supset means "implies" or "if . . . then"), or to argue similarly that

$$(p \gg q \cdot r \gg q) \supset [(p \cdot r) \gg q].$$

For example, "If eating pears (p) is preferable to eating quinces (q) and eating rhubarb (r) is preferable to eating quinces, then eating both pears and rhubarb is preferable to eating quinces." The task is one of erecting a foundation for the systematization of the formal rules governing such a propositional preference relation—a foundation

that can be either axiomatic or linguistic (i.e., in terms of a semantical criterion of acceptability).

One procedure—adapted from the ideas of the Finnish philosopher Georg Henrik von Wright (b. 1916), a prolific contributor to applied logic—is as follows: beginning with a basic set of possible worlds (or states of affairs) w_1, w_2, \dots, w_n , all the propositions to be dealt with are first defined with respect to these by the usual logical connectives (\vee, \cdot, \supset , and so on). Given two elementary propositions p and q, there are just the following possibilities: both are true, p is true and q is false, p is false and q is true, or both are false. Corresponding to each of these possibilities is a possible world; thus,

Approach in terms of possible worlds

$$\begin{aligned} w_1 &= p \cdot q \\ w_2 &= p \cdot \sim q \\ w_3 &= \sim p \cdot q \\ w_4 &= \sim p \cdot \sim q. \end{aligned}$$

The truth of p then amounts to the statement that one of the worlds w_1, w_2 obtains, so that p is equivalent to $w_1 \vee w_2$. Moreover, a given basic preference/indifference ordering among the w_i is assumed. On this basis the following general characterization of propositional preference is stipulated: If delta (δ) is taken to represent any (and thus every) proposition independent of p and q, then p is preferable to q ($p \gg q$), if for every such δ it is the case that every possible world in which p and not-q and δ are the case ($p \cdot \sim q \cdot \delta$) is w-preferable to every possible world in which not-p and q and δ is the case ($\sim p \cdot q \cdot \delta$)—i.e., when $p \cdot \sim q$ is always preferable to $\sim p \cdot q$ provided that everything else is equal. It is readily shown that through this approach such general rules as the following are obtained:

1. If p is preferable to q, then q is not preferable to p; i.e.,

$$p \gg q \vdash \sim (q \gg p).$$

2. If p is preferable to q, and q is preferable to r, then p is preferable to r; i.e.,

$$(p \gg q \cdot q \gg r) \vdash (p \gg r).$$

3. If p is preferable to q, then not-q is preferable to not-p; i.e.,

$$p \gg q \vdash \sim q \gg \sim p.$$

4. If p is preferable to q, then having p and not-q is preferable to having not-p and q; i.e.,

$$p \gg q \vdash (p \cdot \sim q) \gg (\sim p \cdot q).$$

The preceding construction of preference requires only a preference ordering of the possible worlds. If, however, a measure for both probability and desirability (utility) of possible worlds is given, then one can define the corresponding #-value (see below) of an arbitrary proposition p as the probabilistically weighed utility value of all the possible worlds in which the proposition obtains. As an example, p may be the statement "The Franklin Club caters chiefly to business people," and q the statement "The Franklin Club is sports-oriented." It may then be supposed as given that the following values hold:

World	Probability	Desirability
$w_1 = p \cdot q$	1/6	-2
$w_2 = p \cdot \sim q$	2/6	+1
$w_3 = \sim p \cdot q$	2/6	-1
$w_4 = \sim p \cdot \sim q$	1/6	+3

The #-value of a proposition is determined by first multiplying the probability times the desirability of each world in which the proposition is true and then taking the sum of these. For example, the #-value of p is determined as follows: p is true in each of w_1 and w_2 (and only these); the probability times the desirability of w_1 is $1/6 \times (-2)$, and that of w_2 is $2/6 \times (+1)$; thus $\#(p)$ is $1/6 \times (-2) + 2/6 \times (+1) = 0$. (The #-value corresponds to the decision theorists' notion of expected value.) By this procedure it can easily be determined that

$$\begin{aligned} \#(p) &= 0 & \#(\sim p) &= \frac{1}{6} \\ \#(q) &= -\left(\frac{1}{6}\right) & \#(\sim q) &= \frac{2}{6}. \end{aligned}$$

Since both $\#(p) > \#(q)$ and $\#(\sim q) > \#(\sim p)$, one correspondingly obtains both $p \gg q$ and $\sim q \gg \sim p$ in the example at

The "measure of merit"

issue—*i.e.*, “That the Franklin Club should cater chiefly to business people is preferable to its being sports-oriented” and “Its not being sports-oriented is preferable to its not catering chiefly to business people.” (The result is, of course, relative to the given desirability schedule specified for the various possible-world combinations in the above tabulation.)

A more complex mode of preference results, however, if—when some basic utility measure, $\#(x)$, is given—instead of having $p \succ q$ correspond to the condition that $\#(p) > \#(q)$, it is taken to correspond to $\#(p) - \#(\sim p) > \#(q) - \#(\sim q)$. This mode will be governed by characteristic rules, specifically including all those listed above.

The logic of commands. Some scholars have maintained that there cannot be a logic of commands (instructions, orders), inasmuch as there can be no logic in which validity of inference cannot be defined. Validity, however, requires that the concept of truth be applicable (an argument being valid when its conclusion must be true if its premises are true). But, since commands—and for that matter also instructions, requests, and so on—are neither true nor false, it is argued that the concept of validity cannot be applied, so there can be no valid inference in this sphere. This line of thought, however, runs counter to clear intuitions that arise in specific cases, in which one unhesitatingly reasons from commands and sets of commands. If an examination has the instructions “Answer no fewer than three questions! Answer no more than four questions!” one would not hesitate to say that this implies the instruction, “Answer three or four questions!”

This seeming impasse can be broken, in effect, by importing truth into the sphere of commands through the back door: with any command one can associate its termination statement, which, with future-tense reference, asserts it as a fact that what the command orders will be done. Thus, the command “Shut all the windows in the building!” has the termination statement “All the windows in the building will be shut.” In case of a pure command argument—*i.e.*, one that infers a command conclusion from premises that are all commands—validity can be assessed in the light of the validity of the purely assertoric syllogism composed of the corresponding termination statements. Thus the validity of the command argument given above derives from the validity of the inference from the premises “No fewer than three questions will be answered and no more than four questions will be answered” to the conclusion “Three or four questions will be answered.”

The logical issues of pure command inference can be handled in this manner. But what of the mixed cases in which some statement—premise or conclusion—is not a command?

Special case 1. One mixed case is that in which the premises nontrivially include noncommands, but the inferred conclusion is a command. Some writers have endorsed the rule that there is no validity unless the command conclusion is forthcoming from the command premises alone. This, however, invalidates such seemingly acceptable arguments as “Remove all cats from the area; the shed is in the area; so, remove all cats from the shed.” It is more plausible, however, to stipulate the weaker condition that an inference to a command conclusion cannot count as valid unless there is at least one command premise that is essential to the argument. Subject to this restriction, a straightforward application of the above-stated characterization of validity can again be made. This approach validates the above-mentioned command inference via the validity of the assertion inference: “All cats will be removed from the area; the shed is in the area; so, all cats will be removed from the shed.” (The rule under consideration suffices to block the unacceptable argument from the factual premise “All the doors will be shut” to the command conclusion “Shut all the doors.”)

Special case 2. Another mixed case is that in which the premises nontrivially include commands, but the inferred conclusion is an ordinary statement of fact. Some authorities stipulate that no indicative conclusion can be validly drawn from a set of premises which cannot validly be drawn from the indicative among them alone. This rule would seem to be acceptable, though subject to certain

significant provisos: (1) It must be restricted to categorical rather than conditional commands. “If you want to see one of the world’s tallest buildings, look at the Empire State Building” conveys (*inter alia*) the information that “The Empire State Building is one of the world’s tallest buildings.” (2) Exception must be made for those commands that include in their formulation—explicitly or by way of tacit presupposition—reference to a factual datum. “John, give the book to Tom’s brother Jim” yields the fact that Jim is Tom’s brother; and “John, drive your car home” (= “John, you own a car; drive it home”) yields “John owns a car.” With suitable provisos, however, the rule can be maintained to resolve the issues of the special case in view.

Deontic logic. The propositional modalities relating to normative (or valuational) classifications of actions and states of affairs, such as the permitted, the obligatory, the forbidden, or the meritorious, are characterized as deontic modalities (Greek *deontos*, “of that which is binding”) and systematized in deontic logic. Though this subject was first treated as a technical discipline in 1926, its current active development dates from a paper published in 1951 by von Wright. As a highly abstracted branch of logical theory, it leaves to substantive disciplines—such as ethics and law—the concrete questions of what specific acts or states of affairs are to be forbidden, permitted, or the like (just as deductive logic does not meddle with what contingent issues are true but tells only what follows when certain facts or assumptions about the truth are given). It seeks to systematize the abstract, purely conceptual relations between propositions in this sphere, such as the following: if an act is obligatory, then its performance must be permitted and its omission forbidden. In given circumstances, either any act is permitted itself or its omission is permitted.

The systematization and relation to alethic modal logic. In the systematization of deontic logic, the symbols p, q, r, \dots may be taken to range over propositions dealing both with impersonal states of affairs and with the human acts involved in their realization. Certain special deontic operations can then be introduced: $P(p)$ for “It is permitted that p be the case”; $F(p)$ for “It is forbidden that p be the case”; and $O(p)$ for “It is obligatory that p be the case.” In a systematization of deontic logic, it is necessary to take only one of these three operations as primitive (*i.e.*, as an irreducible given), because the others can then be introduced in terms of it. For example, when P alone is taken as primitive (as is done here), the following can be introduced by definition: “It is obligatory that p ” means “It is not permitted that not- p ,” and “It is forbidden that p ” means “It is not permitted that p ”; *i.e.*,

$$O(p) = \sim P(\sim p) \text{ and } F(p) = \sim P(p).$$

The logical grammar of P is presumably to be such that one wants to insist upon the rule:

$$\text{Whenever } \vdash p \supset q, \text{ then } \vdash P(p) \supset P(q).$$

Further, a basic axiom for such an operator as P is

$$\vdash P(p \supset q) \supset (P(p) \supset P(q)),$$

from which it immediately follows that

$$\text{Whenever } \vdash p \supset q, \text{ then } \vdash P(p) \supset P(q).$$

Example: “Since one’s helping Jones, who has been robbed, entails that one help someone who has been robbed, being permitted to help Jones (who has been robbed) entails that one be permitted to help someone who has been robbed.” This yields such principles as “If both p and q are permitted, then p is permitted and q is permitted” and “If p is permitted, then either p or q is permitted”; *i.e.*,

$$\vdash P(p \cdot q) \supset [P(p) \cdot P(q)] \text{ and } \vdash P(p) \supset P(p \vee q).$$

And, once it is postulated that “A p exists that is permitted”—*i.e.*, $\vdash (\exists p)P(p)$ —then the statement that “It is not permitted that both p and not- p ”—*i.e.*, $\sim P(p \cdot \sim p)$ —is also yielded. Moreover, on any adequate theory of P , it is necessary to have such principles as “Either p or not- p is permitted”; *i.e.*, $\vdash P(p \vee \sim p)$.

On the other hand, certain principles must be rejected, such as “If p is permitted and q is permitted, then both

Truth
involve-
ments

Reduction
to
permissi-
bility
terms

Special
cases

p and q taken together are permitted”—i.e., $\neg [P(p) \cdot P(q)] \supset P(p \cdot q)$, in which \neg symbolizes the rejection of a thesis—and that “if either p or q is permitted, then p is permitted”—i.e., $\neg P(p \vee q) \supset P(p)$. The first of these, accepted unqualifiedly, would lead to the untenable result that there can be no permission-indifferent acts—i.e., no acts such that both they and their omission are permitted—since this would then lead to $P(p \cdot \sim p)$. The second thesis would have the unacceptable result of asserting that, when at least one member of a pair of acts is permitted, then both members are permitted.

Analogy with possibility and necessity

In all respects so far considered, deontic logic is wholly analogous to the already well-developed field of alethic modal logic, which deals with statements of the form “It is possible that . . .” (symbolized M), “It is necessary that . . .” (symbolized L), and so on (see above *Modal logic*), with P in the role of possibility (M) and O in that of necessity (L). This parallel, however, does not extend throughout. In alethic logic, the principle that “necessity implies actuality” obviously holds (i.e., $\vdash Lp \supset p$). But its deontic analogue, that “obligation implies actuality” (i.e., $\vdash Op \supset p$), must be rejected, or rather an analogous thesis holds only in the weakened form that “obligation implies permissibility” (i.e., $\vdash Op \supset Pp$). Controversy exists about the relation of deontic to alethic modal logic, principally in the context of Immanuel Kant’s thesis that “ought implies can” (i.e., $\vdash Op \supset Mp$), but also about the theses *ad impossibile nemo obligatur*—“no one is obliged to do the impossible” (i.e., $\vdash \sim Mp \supset \sim Op$)—and “necessity implies permissibility” (i.e., $Lp \supset Pp$). Although this thesis is generally accepted, some scholars want to strengthen the thesis to “necessity implies obligation” (i.e., $\vdash Lp \supset Op$), or, equivalently, to “permissibility implies possibility” (i.e., $\vdash Pp \supset Mp$), with the result that only what is possible can count as permitted, so that the impossible is forbidden. Some would deny that it is wrong (i.e., impermissible) to act to realize the impossible, rather than merely unwise.

It has been proposed that deontic logic may perhaps be reduced to alethic modal logic. This approach is based on the idea of a normative code delimiting the range of the permissible. In this context, what signalizes an action as impermissible is that it involves a violation of the code: the statement that the action has occurred entails that the code has been violated and so leads to a “sanction.” This line of thought leads to the definition of a modal operator $Fp = L(p \supset \sigma)$, “ p necessarily implies a sanction,” in which sigma (σ) is the sanction produced by code violation. Correspondingly, one then obtains “For p to be permitted means that p does not imply by necessity a sanction”—i.e., $Pp = \sim L(p \supset \sigma)$ —and “For p to be obligatory means that not doing p implies by necessity a sanction”—i.e., $Op = L(\sim p \supset \sigma)$. Assuming a systematization of the alethic modal operator L , these definitions immediately produce a corresponding system of deontic logic that—if L is a normal modality—has many of the features that are desirable in a modal operator. It also yields, however—through the “paradoxes of strict implication” (see above *Alternative systems of modal logic*)—the disputed principle that “The assumption that p is not possible implies that p is not permissible”; i.e., $\vdash \sim Mp \supset \sim Pp$. This and other similar consequences of the foregoing effort to reduce deontic logic to modal logic have been transcended by other scholars, who have resorted to a mode of implication (symbolized as \rightarrow) that is stronger than strict implication (as necessary material implication is called) and then defining Fp as $p \rightarrow \sigma$ instead of as above.

Alternative deontic systems. Each of the three principal deontic systems that have been studied to date is analogous to one of the alethic modal systems that were developed in the mid-20th century.

Analogies with M, S4, and S5

These foundational alethic systems differ by virtue of the different axioms and rules adopted for such modalities as necessity, possibility, and contingency. In the system designated M , for example, developed by the aforementioned Finnish logician G.H. von Wright, the adverb “possibly,” symbolized M , is taken as the fundamental undefined modality in terms of which the other modalities are constructed. “Necessarily p ,” symbolized Lp , for example, is defined in the system M as “not possibly not- p ”; i.e.,

$Lp = \sim M \sim p$. Alternatively, in an equivalent system, T , “necessarily p ” is taken as primitive, and “possibly p ” is defined as “not necessarily not- p ”; i.e., $Mp = \sim L \sim p$. Several nonequivalent systems have been developed by the conceptual pragmatist C.I. Lewis (1883–1964), primary author of *Symbolic Logic* (1932), the foundational work in this field. Of these systems, that known as $S4$ includes all of the system M but adds also the axiom that “‘Necessarily p ’ implies ‘It is necessary that necessarily p ’”—i.e., $Lp \supset LLp$ —whereas that known as $S5$ adds still another axiom, that “‘Possibly p ’ implies ‘It is necessary that possibly p ’”—i.e., $Mp \supset LMp$ (see above *Alternative systems of modal logic*). The analogous deontic systems are then as follows:

1. **DM** (the deontic analogue of the system M of von Wright or of the system T). To a standard system of propositional logic the following rule is added: “Any proposition, if true, ought to be true”; that is, $\vdash p$ then $\vdash Op$. Example: Given that “to forgive is divine” (p), then “to forgive ought to be divine” (Op). Axioms:
 - A1. “If p is obligatory, then not- p is not obligatory”; i.e., $Op \supset \sim O \sim p$.
 - A2. “If p ought to imply q , then if p is obligatory q is obligatory”; i.e., $O(p \supset q) \supset (Op \supset Oq)$.
2. **DS4** (the deontic analogue of Lewis’ system $S4$). To M one adds the axiom:
 - A3. “If p is obligatory, then p ought to be obligatory”; i.e., $Op \supset OOp$. Example: “If John ought to pay his debts” (Op), then it is obligatory that John ought to pay his debts” (OOp).
3. **DS5** (the deontic analogue of Lewis’ system $S5$). To M one adds the axiom:
 - A4. “If p is not obligatory, then p ought to be nonobligatory”; i.e., $\sim Op \supset O \sim Op$.

Semantic systematization

A straightforward semantical systematization of systems of deontic logic can be provided as follows: given a domain of complex propositions built up from atomic propositions (p, q, r, \dots) with the use of propositional connectives ($\sim, \cdot, \vee, \supset$) and O , a deontic model set Δ for this domain can be characterized as any set chosen from these propositions that meets the following conditions (in which “iff” means “if and only if”):

1. Not- p is in the set if and only if p is not in the set; i.e., $\sim p \in \Delta$ iff $p \notin \Delta$.
2. “Both p and q together” is in the set if and only if p is in the set and q is in the set; i.e., $(p \cdot q) \in \Delta$ iff $p \in \Delta$ and $q \in \Delta$.
3. “Either p or q ” is in the set if and only if either p is in the set or q is in the set; i.e., $(p \vee q) \in \Delta$ iff $p \in \Delta$ or $q \in \Delta$.
4. “That p implies q ” is in the set if and only if either p is not in the set or q is in the set; i.e., $(p \supset q) \in \Delta$ iff $p \notin \Delta$ or $q \in \Delta$.
5. “That p is obligatory” is in the set whenever p is posited; i.e., $Op \in \Delta$ whenever $\vdash p$.
6. “That not- p is not obligatory” is in the set whenever “ p is obligatory” is in the set; i.e., $\sim O \sim p \in \Delta$ whenever $Op \in \Delta$.
7. “That q is obligatory” is in the set whenever both “ p is obligatory” is in the set and “That p implies q is obligatory” is in the set; i.e., $Oq \in \Delta$ whenever both $Op \in \Delta$ and $O(p \supset q) \in \Delta$.

A proposition can be characterized as a deontic thesis (D-thesis) if it can be shown that, in virtue of these rules, it must belong to every deontic model set. It can be demonstrated that the D-thesis in this sense will coincide exactly with the theorems of **DM**—the first of the above three systems. Furthermore, if one adds one of the additional rules:

- 8’. “That p ought to be obligatory” is in the set whenever “ p is obligatory” is in the set; i.e., $OOp \in \Delta$ whenever $Op \in \Delta$.
- 8”. “That p ought to be non-obligatory” is in the set whenever “ p is not obligatory” is in the set; i.e., $O \sim Op \in \Delta$ whenever $\sim Op \in \Delta$.

then the corresponding D' or D'' theses will coincide exactly with the theorems of the deontic systems **DS4** and **DS5**, respectively—numbers 2 and 3 above.

LOGICS OF PHYSICAL APPLICATION

Certain systems of logic are built up specifically with particular physical applications in view. Within this range lie temporal logic; spatial, or topological, logic; mereology, or the logic of parts and wholes generally; as well as the logic of circuit analysis.

Since the field of topological logic is still relatively undeveloped, the reader is referred to the bibliography for a recent source that provides some materials and references to the literature.

Temporal logic. The object of temporal logic—variously called chronological logic or tense logic—is to systematize reasoning with time-related propositions. Such propositions generally do not involve the timeless “is” (or “are”) of the mathematicians’ “three is a prime,” but rather envisage an explicitly temporal condition (examples: “Bob is sitting,” “Robert was present,” “Mary will have been informed”). In this area, statements are employed in which some essential reference to the before-after relationship or the past-present-future relationship is at issue; and the ideas of succession, change, and constancy enter in.

Classic historical treatments. Chronological logic originated with the Megarians of the 4th century BC, whose school (not far from Athens) reflected the influence of Socrates and of Eleaticism.

In the Megarian conception of modality, the actual is that which is realized now, the possible is that which is realized at some time or other, and the necessary is that which is realized at all times. These Megarian ideas can be found also in Aristotle, together with another temporalized sense of necessity according to which certain possibilities are possible prior to the event, actual then, and necessary thereafter, so that their modal status is not omnitemporal (as in the Megarian concept) but changes in time. The Stoic conception of temporal modality is yet another cognate development, according to which the possible is that which is realized at some time in the present or future, and the necessary that which is realized at all such times. The Diodorean concept of implication (named after the 4th-century-BC Megarian logician Diodorus Cronus) holds, for example, that the conditional “If the sun has risen, it is daytime” is to be given the temporal construction “All times after the sun has risen are times when it is daytime.” The Persian logician Avicenna (980–1037), the foremost philosopher of medieval Islām, treated this chronological conception of implication in the framework of a general theory of categorical propositions (such as “All A is B”) of a temporalized type and considerably advanced and developed the Megarian-Stoic theory of temporal modalities.

Fundamental concepts and relations of temporal logic. The statements “It sometimes rains in London,” “It always rains in London,” and “It is raining in London on Jan. 1, AD 3000,” are all termed chronologically definite, in that their truth or falsity is independent of their time assertion. By contrast, the statements “It is now raining in London,” “It rained in London yesterday,” and “It will rain in London sometime next week” are all chronologically indefinite, in that their truth or falsity is not independent of their time of assertion. The notation $| t \vdash p$ is here introduced to mean that the proposition p , often in itself chronologically indefinite, is represented as being asserted at the time t . For example, if p_1 is the statement “It is raining in London today” and t_1 is Jan. 1, 1900, then “ $| t_1 \vdash p_1$ ” represents the assertion made on Jan. 1, 1900, that it is raining today—an assertion that is true if and only if the statement “It is raining in London on Jan. 1, 1900,” is true. If the statement p is chronologically definite, then (by definition) the assertions “ $| t \vdash p$ ” and “ $| t' \vdash p$ ” are materially equivalent (*i.e.*, have the same truth value) for all values of t and t' . Otherwise, p is chronologically indefinite. The time may be measured, for example, in units of days, so that the time variable is made discrete. Then $(t + 1)$ will represent “the day after t -day,” $(t - 1)$ will represent “the day before t -day,” and the like. And, further, the statements p_1 , q_1 , and r_1 can then be as follows:

- p_1 : “It rains in London today.”
- q_1 : “It will rain in London tomorrow.”
- r_1 : “It rained in London yesterday.”

The following assertions can now be made:

- P: $| t \vdash p_1$
- Q: $| t - 1 \vdash q_1$
- R: $| t + 1 \vdash r_1$.

Clearly, for any value of t whatsoever, the assertions P, Q, and R must (logically) be materially equivalent (*i.e.*, have the same truth value). This illustration establishes the basic point—that the theory of chronological propositions must be prepared to exhibit the existence of logical relationships among propositions of such a kind that the truth of the assertion of one statement at one time may be bound up essentially with the truth (or falsity) of the assertion of some very different statement at another time.

A (genuine) date is a time specification that is chronologically stable (such as “Jan. 1, 3000,” or “the day of Lincoln’s assassination”); a pseudodate is a time specification that is chronologically unstable (such as “today” or “six weeks ago”). These lead to very different results depending on the nature of the fundamental reference point—the “origin” in mathematical terms. If the origin is a pseudodate—say, “today”—the style of dating will be such that its chronological specifiers are pseudodates—tomorrow, the day before yesterday, four days ago, and so on. If, on the other hand, the origin is a genuine date, say that of the founding of Rome or the accession of Alexander, the style of dating will be such that all its dates are of the type: two hundred and fifty years *ab urbe condita* (“since the founding of the city”). Clearly, a chronology of genuine dates will then be chronologically definite, and one of pseudodates will be chronologically indefinite.

Let p be some chronologically indefinite statement. Then, in general, another statement can be formed, asserting that p holds (obtains) at the time t . Correspondingly, let the statement-forming operation R_t be introduced. The statement $R_t(p)$, which is to be read “ p is realized at the time t ,” will then represent the statement stating explicitly that p holds (obtains) specifically at the time t . Thus, if t_1 is 3:00 PM Greenwich Mean Time on Jan. 1, 2000, and p_1 is the (chronologically indefinite) statement “All men are (*i.e.*, are now) playing chess,” then “ $R_{t_1}(p_1)$ ” is the statement “It is the case at 3:00 PM Greenwich Mean Time on Jan. 1, 2000, that all men are playing chess.”

Systematization of temporal reasoning. On the basis of these ideas, the logical theory of chronological propositions can be developed in a systematic, formal way. It may be postulated that the operator R is to be governed by the following rules:

The negation of a statement p is realized at a given time if and only if it is not the case that the statement is realized at that time; *i.e.*, $R_t(\neg p) \equiv \neg R_t(p)$, in which \equiv signifies equivalence and is read “if and only if.” (T1)

A conjunction of two statements is realized at a given time if and only if each of these two statements is realized at that time: $R_t(p \cdot q) \equiv [R_t(p) \cdot R_t(q)]$. Example: “John and Jane are at the railroad station at 10:00 AM— $R_t(p \cdot q)$ —if and only if John is at the station at 10:00 AM— $R_t(p)$ —and Jane is at the station at 10:00 AM— $R_t(q)$.” (T2)

If a statement is realized universally—*i.e.*, at any and every time whatsoever—it can then be expressed more simply as being true without any temporal qualifications; hence the rule:

If for every time t the statement p is realized, then p obtains unqualifiedly; *i.e.*, $(\forall t)R_t(p) \supset p$, in which \forall is the universal quantifier.

If two times are involved, however, then the left-hand term in rule (T3) can be expressed within the second time frame as “It will be the case τ from now that, for every time t , it will be the case t from the first now that p ”; *i.e.*, $R_t[(\forall t)R_t(p)]$. It is an algebraic rule, however, that an R_t operator can be moved to the right past an irrelevant quantifier; hence

Dated versus pseudodated statements

Definite versus indefinite reference

$$R_t[(\forall t)R_t(p)] \equiv (\forall t)(R_t[R_t(p)]);$$

and, correspondingly, with the existential quantifier \exists : "It will be the case τ from now that there exists a time t such that p will be realized at t " is equivalent to saying "There exists a time t such that it will be the case τ from now that p will be realized t from the first now" (in which τ is a second time); *i.e.*,

$$R_t[(\exists t)R_t(p)] \equiv (\exists t)(R_t[R_t(p)]). \quad (T4)$$

It is notable that the left-hand side of this equivalence is itself equivalent with $(\exists t)R_t(p)$ since what follows the initial R_t is a chronologically definite statement.

Finally, there are two distinct ways of construing iterations of the R_t operator, depending on the choice of origin of the second time scale. Thus a choice is required between two possible rules:

$$R_t[R_t(p)] \equiv R_t(p) \quad (T5-I)$$

$$R_t[R_t(p)] \equiv R_{t+\Delta}(p). \quad (T5-II)$$

Taking these rules as a starting point, two alternative axiomatic theories are generated for the logic of the operation of chronological realization.

Rules and modalities

Apart from strictly technical results establishing the formal relationships between the various systems of chronological logic, the most interesting findings about the systems of tense logic relate to the theory of temporal modalities. The most striking finding concerns the logical structure of the system of modalities, be it Megarian or Stoic:

Megarian	{possibly p : $(\exists t)R_t(p)$
	{necessarily p : $(\forall t)R_t(p)$
Stoic	{possibly p : $(\exists t)[F(t) \cdot R_t(p)]$
	{necessarily p : $(\forall t)[F(t) \supset R_t(p)]$

in which $F(t)$ signifies "t is future." It has been shown that the forms, or structures, of both of these systems of temporal modalities are given by the aforementioned system $S5$ of C.I. Lewis. Exactly parallel results are obtained for modalities of past times, $P_t(p)$: p was realized at some (past) time t ; and $\sim P_t(\sim p)$: p has been realized at all (past) times.

Mereology. The founder of mereology was the Polish logician Stanisław Leśniewski. Leśniewski was much exercised about Russell's paradox of the class of all classes not elements of themselves—if this class is a member of itself, then it is not; and if it is not, then it is (example: "This barber shaves everyone in town who does not shave himself." Does the barber then shave himself? If he does, he does not; if he does not, he does.).

Basic concepts and definitions. The paradox results, Leśniewski argued, from a failure to distinguish the distributive and the collective interpretations of class expressions. The statement " x is an element of the class of X 's" is correspondingly equivocal. When its key terms (element of, class of) are used distributively, it means simply that x is an X . But, if these terms are used collectively, it means that x is a part (proper or improper) of the whole consisting of the X 's—*i.e.*, that x is a part of the object that meets the following two conditions: (1) that every x is a part of it and (2) that every part of it has a common part with some x . On either construction of class membership, one of the inferences essential to the derivation of Russell's paradox is blocked.

Leśniewski presented his theory of the collective interpretation of class expressions in a paper published in 1916. Eschewing symbolization, he formulated his theorems and their proofs in ordinary language. Later he sought to formalize the theory by embedding it within a broader body of logical theory. This theory comprised two parts: protothetic, a logic of propositions (not analyzed into their parts); and ontology, which contains counterparts to the predicational logic (of subjects and predicates), including the calculus of relations and the theory of identity. On his own approach, mereology was developed as an extension of ontology and protothetic, but the practice of most later writers has been to develop as a counterpart to mereology a theory of parts and wholes that is simply an extension of the more familiar machinery of quantificational logic employing \exists and \forall . This is the course adopted here.

An undefined relation Pt serves as the basis for an axiomatic theory of the part relation. This relation is operative with respect to the items of some domain D , over which the variables $\alpha, \beta, \gamma, \dots$ (alpha, beta, gamma, and so on) are assumed to range. Thus, $\alpha Pt \beta$ is to be read "alpha is a part of beta"—with "part" taken in the wider sense in which the whole counts as part of itself. Two definitions are basic:

Disjointness and summation

1. " α is disjoint from β "; *i.e.*, $\alpha | \beta$ is defined as obtaining when "there exists no item γ such that γ is a part of α and γ is a part of β "; *i.e.*, $\sim(\exists \gamma)(\gamma Pt \alpha \cdot \gamma Pt \beta)$. Example: "The transmission (α) is disjoint from the motor (β) if there exists no machine part (γ) such that it is a part of the transmission and also a part of the motor."
2. " S has the sum of (or sums to) α "; *i.e.*, $S \Sigma \alpha$ is defined as obtaining when "for every γ , this γ is disjoint from α if and only if, for every β , to be a member of S is to be disjoint from γ "; *i.e.*,

$$(\forall \gamma)[\gamma | \alpha \equiv (\forall \beta)(\beta \in S \supset \beta | \gamma)].$$

$S \Sigma \alpha$ thus obtains whenever everything disjoint from α is disjoint from every S -element (β) as well, and conversely. Example: "A given group of buildings (S) comprises (Σ) the University of Oxford (α) when, for every room in the world (γ)—office, classroom, etc.—this room is disjoint from the university if and only if, in the case of each building (β), for it to be a member (\in) of the group that comprises the university (S) it must not have this room as a part ($\beta | \gamma$)."

Axiomatization of mereology. A comprehensive theory of parts and wholes can now be built up from three axioms:

The first axiom expresses the fact that "for every α and every β , if α is a part of β and β is a part of α , then α and β must be one and the same item"; *i.e.*,

$$(\forall \alpha)(\forall \beta)(\alpha Pt \beta \cdot \beta Pt \alpha \supset \alpha = \beta);$$

hence, the axiom:

$$\text{Items that are parts of one another are identical.} \quad (A1)$$

The second axiom expresses the fact that "for every α and every β , α is a part of β if and only if, for every γ , if this γ is disjoint from β it is then disjoint from α as well"; *i.e.*,

$$(\forall \alpha)(\forall \beta)[\alpha Pt \beta \equiv (\forall \gamma)(\gamma | \beta \supset \gamma | \alpha)];$$

hence, the axiom:

$$\text{One item is part of another only if every item disjoint from the second is also disjoint from the first.} \quad (A2)$$

The third axiom expresses the fact that "if there exists an α that is a member of a nonempty set of items S , then there also exists a β that is the sum of this set"; *i.e.*,

$$(\exists \alpha)(\alpha \in S) \supset (\exists \beta)S \Sigma \beta;$$

hence, the axiom:

$$\text{Every nonempty set has a sum.} \quad (A3)$$

Several theorems follow from these axioms:

Deduced theorems

The first states that "for every α , α is a part of α "; *i.e.*,

$$(\forall \alpha)\alpha Pt \alpha;$$

hence, the theorem:

$$\text{Every item is part of itself.} \quad (T1)$$

The second theorem states that "for every α , for every β , and for every γ , if α is a part of β , and β is a part of γ , then α is a part of γ "; *i.e.*,

$$(\forall \alpha)(\forall \beta)(\forall \gamma)[(\alpha Pt \beta \cdot \beta Pt \gamma) \supset \alpha Pt \gamma];$$

hence, the theorem:

$$\text{The Pt-relation is transitive.} \quad (T2)$$

The third theorem states that "for every α , for every β , and for every γ , if γ is a part of α only when it is also a part of β , then α is identical with β "; *i.e.*,

$$(\forall \alpha)(\forall \beta)(\forall \gamma)[(\gamma Pt \alpha \equiv \gamma Pt \beta) \supset \alpha = \beta];$$

hence, the theorem:

Any item is completely determined by its parts; items are identical when they have the same parts in common. (T3)

The fourth theorem states that "for every α and every β , there exists a γ that is the sum of α and β "; i.e.,

$$(\forall\alpha)(\forall\beta)(\exists\gamma)((\alpha, \beta)\Sigma\gamma);$$

hence, the theorem:

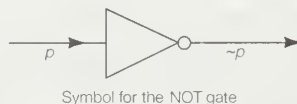
Any two items whatsoever may be summed up. (T4)

In this form as a formal theory of the part relation, the history of mereology can be dated from some drafts and essays of Leibniz prepared in the late 1690s.

COMPUTER DESIGN AND PROGRAMMING

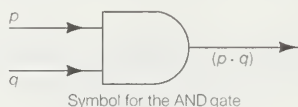
In the most general terms a computer is a device that calculates a result ("output") from one or more initial items of information ("input"). Inputs and outputs are usually represented in binary terms—i.e., in strings of 0s and 1s—and the values of 0 and 1 are realized in the machine by the presence or absence of a current (of electricity, water, light, and so on). When the output is a completely determined function of the input, the connection between a computer and the two-valued logic of propositions is immediate, for a valid argument can be construed as a partial function of the truth values of the premises such that when the premises each have the value true, so does the conclusion.

One of the simplest computers has one input, either 0 or 1 (i.e., a current either off or on), and one output, namely, the reverse of the input. That is, when 0 is input, 1 is output, and, conversely, when 1 is input, 0 is output. This is also the behaviour of the truth function negation ($\sim p$) when applied to the truth values true and false. Thus a circuit element that behaves in such a way is called a NOT gate:

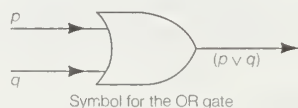


When no current is input from the left, a current flows out on the right, and, conversely, when a current flows in from the left, none is output to the right.

Similarly, devices with two inputs and one output correspond in behaviour to the truth functions conjunction ($p \cdot q$) and disjunction ($p \vee q$). Specifically, in an AND gate,



current flows out to the right only when current is present in both inputs; otherwise there is no output. In an OR gate, current is output when a current is present in either or both of the inputs on the left.



Other truth functional connectives are easily constructed using combinations of these gates. For example, the conditional, ($p \supset q$), is represented by:



There is no output if there is input from p (" p " is true) and none from q (" q " is false).

It is also possible to connect these gates to memory devices that store intermediate results in order to construct circuits that perform elementary binary arithmetic: addition, subtraction, multiplication, and division. These

simple circuits, and others like them, can be connected together in order to perform various computations such as determining the implications of a set of premises or determining the numerical value of a mathematical function for specific argument values.

The details of computer design and architecture depend less on logical theory and more on the mathematical theory of lattices (see ALGEBRA: *Lattice theory*) and are outside the scope of this article. In computer programming, however, logic has a significant role.

Some modern computers, such as the ones in automobiles or washing machines, are dedicated; that is, they are constructed to perform only certain sorts of computations. Others are general-purpose computers, which require a set of instructions about what to do and when to do it. A set of such instructions is called a program. A general-purpose computer operating under a program begins in an initial state with a given input, passes through intermediate states, and should eventually stop in a final state with a definite output. For a given program, the various momentary states of the machine are characterized by the momentary values of all the variables in the program.

In 1974 the British computer scientist Rod M. Burstall first remarked on the connection between machine states and the possible worlds used in the semantics of modal logic (see above *Modal logic*). The use of concepts and results from modal logic to investigate the properties and behaviour of computer programs (e.g., does this program stop after a finite number of steps?) was soon taken up by others, notably Vaughan R. Pratt (dynamic logic), Amir Pnueli (temporal logic), and David Harel (process logic).

The connection between the possible worlds of the logician and the internal states of a computer is easily described. In possible world semantics, p is possible in some world w if and only if p is true in some world w' accessible to w . Depending on the properties of the accessibility relation (reflexive, symmetric, and so on), there will be different theorems about possibility and necessity (" p is necessary" = " $\sim M \sim p$ "). The accessibility relation of modal logic semantics can thus be understood as the relation between states of a computer under the control of a program such that, beginning in one state, the machine will (in a finite time) be in one of the accessible states. In some programs, for instance, one cannot return from one state to an earlier state: hence state accessibility here is not symmetric. (For detailed treatments of this subject, refer to the *Bibliography*.)

HYPOTHETICAL REASONING AND COUNTERFACTUAL CONDITIONALS

A simple conditional, or "if," statement asserts a strictly formal relationship between antecedent ("if" clause) and consequent ("then" clause): "If p , then q ," without any reference to the status of the antecedent. The knowledge status of this antecedent, however, may be problematic (unknown), or known-to-be-true, or known-to-be-false. In these three cases, one obtains, respectively, the problematic conditional ("Should it be the case that p —which it may or may not be—then q "), the factual conditional ("Since p , then q "), and the counterfactual conditional ("If it were the case that p —which it is not—then q "). Counterfactual conditionals have a special importance in the area of thought experiments in history as well as elsewhere.

Material implication, $p \supset q$, construed simply as the truth-functional "either not- p or q ," is clearly not suited to represent counterfactual conditionals, because any material implication with a false antecedent is true: when p is false, then $p \supset q$ and $p \supset \sim q$ are both true, regardless of what one may choose to put in place of q . But even when a stronger mode of implication is invoked, such as strict implication or its cognates, the problem of auxiliary hypotheses (soon to be explained) would still remain.

It seems most natural to view a counterfactual conditional in the light of an inference to be drawn from the contrary-to-fact thesis represented by its antecedent. Thus, "If this rubber band were made of copper, then it would conduct electricity" would be construed as an incomplete presentation of the argument resulting from its expansion into:

Assumption: "This rubber band is made of copper."
 Known fact: "Everything made of copper conducts electricity."
 Conclusion: "This rubber band conducts electricity."

On the analysis, the conclusion (= the consequent of the counterfactual) appears as a deductive consequence of the assumption (= the antecedent of the counterfactual). This truncated-argument analysis of counterfactuals is a contribution, in essence, of a Polish linguistic theorist, Henry Hiz (b. 1917). On Hiz's analysis, counterfactual conditionals are properly to be understood as metalinguistic—i.e., as making statements about statements. Specifically, "If *A* were so, then *B* would be so" is to be construed in the context of a given system of statements *S*, saying that when *A* is adjoined as a supplemental premise to *S*, then *B* follows. This approach has been endorsed by the American Roderick Chisholm, an important writer in applied logic, and has been put forward by many logicians, most of whom incline to take *S*, as above, to include all or part of the corpus of scientific laws.

Hiz's analysis

The approach warrants a closer scrutiny. On fuller analysis, the following situation, with a considerably enlarged group of auxiliary hypotheses, comes into focus:

- Known facts: 1. "This band is made of rubber."
 2. "This band is not made of copper."
 3. "This band does not conduct electricity."
 4. "Things made of rubber do not conduct electricity."
 5. "Things made of copper do conduct electricity."
 Assumption: Not-2; i.e., "This band is made of copper."

Evaluation through hypotheses, facts, and laws

When this assumption is introduced within the framework of known facts, a contradiction obviously ensues. How

can this situation be repaired? Clearly, the logician must begin by dropping items 1 and 2 and replacing them with their negations—the assumption itself so instructs him. But a contradiction still remains. The following alternatives are open:

- Alternative 1: Retain: 3, 4. Reject: 1, 2, 5.
 Alternative 2: Retain: 4, 5. Reject: 1, 2, 3.

That is, the analyst actually has a choice between rejecting 3 in favour of 5 or 5 in favour of 3, resulting in the following conditionals:

1. "If this rubber band were made of copper, then it would conduct electricity" (since copper conducts electricity).
2. "If this rubber band were made of copper, then copper would not (always) conduct electricity" (since this band does not conduct electricity).

If the first conditional seems more natural than the second, this is owing to the fact that, in the face of the counterfactual hypothesis at issue, the first invites the sacrifice of a particular fact (that the band does not conduct electricity) in favour of a general law (that copper conducts electricity), whereas the second counterfactual would have sacrificed a law to a purely hypothetical fact. On this view, there is a fundamental epistemological difference between actual and hypothetical situations: in actual cases one makes laws give way to facts, but in hypothetical cases one makes the facts yield to laws.

But in more complex cases the fact/law distinction may not help matters. For example, assume a group of three laws L_1, L_2, L_3 , where $\sim L_1$ is inconsistent with the conjunction of L_2 and L_3 . If asked to hypothesize the denial of L_1 —so that the "fact" that one is opposing is itself a law—then what remains is a choice between laws; the distinction between facts and laws does not resolve the issue, and some more sophisticated mechanism for a preferential choice among laws is necessary. (N.R./M.L.Sc.)

THE HISTORY OF LOGIC

Origins of logic in the West

PRECURSORS OF ANCIENT LOGIC

There was a medieval tradition according to which the Greek philosopher Parmenides (5th century BC) invented logic while living on a rock in Egypt. The story is pure legend, but it does reflect the fact that Parmenides was the first philosopher to use an extended argument for his views, rather than merely proposing a vision of reality. But using arguments is not the same as studying them, and Parmenides never systematically formulated or studied principles of argumentation in their own right. Indeed, there is no evidence that he was even aware of the implicit rules of inference used in presenting his doctrine.

Perhaps Parmenides' use of argument was inspired by the practice of early Greek mathematics among the Pythagoreans. Thus it is significant that Parmenides is reported to have had a Pythagorean teacher. But the history of Pythagoreanism in this early period is shrouded in mystery, and it is hard to separate fact from legend.

If Parmenides was not aware of general rules underlying his arguments, the same perhaps is not true for his disciple Zeno of Elea (5th century BC). Zeno was the author of many arguments, known collectively as "Zeno's Paradoxes," purporting to infer impossible consequences from a non-Parmenidean view of things and so to refute such a view and indirectly to establish Parmenides' monist position. The logical strategy of establishing a claim by showing that its opposite leads to absurd consequences is known as *reductio ad absurdum*. The fact that Zeno's arguments were all of this form suggests that he recognized and reflected on the general pattern.

Zeno's Paradoxes

Other authors too contributed to a growing Greek interest in inference and proof. Early rhetoricians and sophists—e.g., Gorgias, Hippias, Prodicus, and Protagoras (all 5th-century BC)—cultivated the art of defending or attacking a thesis by means of argument. This concern for the

techniques of argument on occasion merely led to verbal displays of debating skills, what Plato called "eristic." But it is also true that the sophists were instrumental in bringing argumentation to the central position it came uniquely to hold in Greek thought. The sophists were, for example, among the first people anywhere to demand that moral claims be justified by reasons.

Certain particular teachings of the sophists and rhetoricians are significant for the early history of logic. For example, Protagoras is reported to have been the first to distinguish different kinds of sentences: questions, answers, prayers, and injunctions. Prodicus appears to have maintained that no two words can mean exactly the same thing. Accordingly, he devoted much attention to carefully distinguishing and defining the meanings of apparent synonyms, including many ethical terms.

Socrates (c. 470–399 BC) is said to have attended Prodicus' lectures. Like Prodicus, he pursued the definitions of things, particularly in the realm of ethics and values. These investigations, conducted by means of debate and argument as portrayed in the writings of Plato (428/427–348/347 BC), reinforced Greek interest in argumentation and emphasized the importance of care and rigour in the use of language.

Plato continued the work begun by the sophists and by Socrates. In the *Sophist*, he distinguished affirmation from negation and made the important distinction between verbs and names (including both nouns and adjectives). He remarked that a complete statement (*logos*) cannot consist of either a name or a verb alone but requires at least one of each. This observation indicates that the analysis of language had developed to the point of investigating the internal structures of statements, in addition to the relations of statements as a whole to one another. This new development would be raised to a high art by Plato's pupil Aristotle (384–322 BC).

There are passages in Plato's writings where he suggests

that the practice of argument in the form of dialogue (Platonic "dialectic") has a larger significance beyond its occasional use to investigate a particular problem. The suggestion is that dialectic is a science in its own right, or perhaps a general method for arriving at scientific conclusions in other fields. These seminal but inconclusive remarks indicate a new level of generality in Greek speculation about reasoning.

ARISTOTLE

The logical work of all these men, important as it was, must be regarded as piecemeal and fragmentary. None of them was engaged in the systematic, sustained investigation of inference in its own right. That seems to have been done first by Aristotle. At the end of his *Sophistic Refutations*, Aristotle acknowledges that in most cases new discoveries rely on previous labours by others, so that, while those others' achievements may be small, they are seminal. But then he adds:

Of the present inquiry, on the other hand, it was not the case that part of the work had been thoroughly done before, while part had not. Nothing existed at all. . . . [O]n the subject of deduction we had absolutely nothing else of an earlier date to mention, but were kept at work for a long time in experimental researches.

(From *The Complete Works of Aristotle: The Revised Oxford Translation*, ed. Jonathan Barnes, 1984, by permission of Oxford University Press.)

The
Organon

Aristotle's logical writings comprise six works, known collectively as the *Organon* ("Tool"). The significance of the name is that logic, for Aristotle, was not one of the theoretical sciences. These were physics, mathematics, and metaphysics. Instead, logic was a tool used by all the sciences. (To say that logic is not a science in this sense is in no way to deny it is a rigorous discipline. The notion of a science was a very special one for Aristotle, most fully developed in his *Posterior Analytics*.)

Aristotle's logical works, in their traditional but not chronological order, are:

1. *Categories*, which discusses Aristotle's 10 basic kinds of entities: substance, quantity, quality, relation, place, time, position, state, action, and passion. Although the *Categories* is always included in the *Organon*, it has little to do with logic in the modern sense.
2. *De interpretatione* (*On Interpretation*), which includes a statement of Aristotle's semantics, along with a study of the structure of certain basic kinds of propositions and their interrelations.
3. *Prior Analytics* (two books), containing the theory of syllogistic (described below).
4. *Posterior Analytics* (two books), presenting Aristotle's theory of "scientific demonstration" in his special sense. This is Aristotle's account of the philosophy of science or scientific methodology.
5. *Topics* (eight books), an early work, which contains a study of nondemonstrative reasoning. It is a miscellany of how to conduct a good argument.
6. *Sophistic Refutations*, a discussion of various kinds of fallacies. It was originally intended as a ninth book of the *Topics*.

Aristotle's logic was a term logic, in the following sense. Consider the schema: "If every β is an α and every γ is a β , then every γ is an α ." The " α ," " β ," and " γ " are variables—i.e., placeholders. Any argument that fits this pattern is a valid syllogism and, in fact, a syllogism in the form known as Barbara. (On this terminology, see below.)

The variables here serve as placeholders for terms or names. Thus, replacing " α " by "substance," " β " by "animal," and " γ " by "dog" in the schema yields: "If every animal is a substance and every dog is an animal, then every dog is a substance," a syllogism in Barbara. Aristotle's logic was a term logic in the sense that it focused on logical relations among such terms in valid inferences.

Aristotle was the first logician to use variables. This innovation was tremendously important, since without them it would have been impossible for him to reach the level of generality and abstraction that he did.

Most of Aristotle's logic was concerned with certain kinds of propositions that can be analyzed as consisting of

(1) usually a quantifier ("every," "some," or the universal negative quantifier "no"), (2) a subject, (3) a copula, (4) perhaps a negation ("not"), (5) a predicate. Propositions analyzable in this way were later called categorical propositions and fall into one or another of the following forms:

1. Universal affirmative: "Every β is an α ."
2. Universal negative: "Every β is not an α ," or equivalently "No β is an α ."
3. Particular affirmative: "Some β is an α ."
4. Particular negative: "Some β is not an α ."
5. Indefinite affirmative: " β is an α ."
6. Indefinite negative: " β is not an α ."
7. Singular affirmative: " x is an α ," where " x " refers to only one individual (e.g., "Socrates is an animal").
8. Singular negative: " x is not an α ," with " x " as before.

Sometimes, and very often in the *Prior Analytics*, Aristotle adopted alternative but equivalent formulations. Instead of saying, for example, "Every β is an α ," he would say, " α belongs to every β " or " α is predicated of every β ."

In syllogistic, singular propositions (affirmative or negative) were generally ignored, and indefinite affirmatives and negatives were treated as equivalent to the corresponding particular affirmatives and negatives. In the Middle Ages, propositions of types 1–4 were said to be of forms A, E, I, and O, respectively. This notation will be used below.

In the *De interpretatione* Aristotle discussed ways in which affirmative and negative propositions with the same subjects and predicates can be opposed to one another. He observed that when two such propositions are related as forms A and E, they cannot be true together but can be false together. Such pairs Aristotle called contraries. When the two propositions are related as forms A and O or as forms E and I or as affirmative and negative singular propositions, then it must be that one is true and the other false. These Aristotle called contradictories. He had no special term for pairs related as forms I and O, although they were later called subcontraries. Subcontraries cannot be false together, although, as Aristotle remarked, they may be true together. The same holds for indefinite affirmatives and negatives, construed as equivalent to the corresponding particular forms. Note that if a universal proposition (affirmative or negative) is true, its contradictory is false, and so the subcontrary of that contradictory is true. Thus propositions of form A imply the corresponding propositions of form I, and those of form E imply those of form O. These last relations were later called subalternation, and the particular propositions (affirmative or negative) were said to be subalternate to the corresponding universal propositions.

Near the beginning of the *Prior Analytics*, Aristotle formulated several rules later known collectively as the theory of conversion. To "convert" a proposition in this sense is to interchange its subject and predicate. Aristotle observed that propositions of forms E and I can be validly converted in this way: if no β is an α , then so too no α is a β , and if some β is an α , then so too some α is a β . In later terminology, such propositions were said to be converted "simply" (*simpliciter*). But propositions of form A cannot be converted in this way: if every β is an α , it does not follow that every α is a β . It does follow, however, that some α is a β . Such propositions, which can be converted provided that not only are their subjects and predicates interchanged but also the universal quantifier is weakened to a particular quantifier "some," were later said to be converted "accidentally" (*per accidens*). Propositions of form O cannot be converted at all: from the fact that some animal is not a dog, it does not follow that some dog is not an animal. Aristotle used these laws of conversion in later chapters of the *Prior Analytics* to reduce other syllogisms to syllogisms in the first figure, as described below.

Aristotle defined a syllogism as "discourse in which, certain things being stated something other than what is stated follows of necessity from their being so." (From *The Complete Works of Aristotle: The Revised Oxford Translation*, ed. Jonathan Barnes, 1984, by permission of Oxford University Press.) But in practice he confined the term to arguments containing two premises and a conclusion, each of which is a categorical proposition. The

Categorical
forms

Theory
of
conversion

subject and predicate of the conclusion each occur in one of the premises, together with a third term (the middle) that is found in both premises but not in the conclusion. A syllogism thus argues that because α and γ are related in certain ways to β (the middle) in the premises, they are related in a certain way to one another in the conclusion.

The predicate of the conclusion is called the major term, and the premise in which it occurs is called the major premise. The subject of the conclusion is called the minor term and the premise in which it occurs is called the minor premise. This way of describing major and minor terms conforms to Aristotle's actual practice and was proposed as a definition by the 6th-century Greek commentator John Philoponus. But in one passage Aristotle put it differently: the minor term is said to be "included" in the middle and the middle "included" in the major term. This remark, which appears to have been intended to apply only to the first figure (see below), has caused much confusion among some of Aristotle's commentators, who interpreted it as applying to all three figures.

Aristotle distinguished three different figures of syllogisms, according to how the middle is related to the other two terms in the premises. In one passage, he says that if one wants to prove α of γ syllogistically, one finds a middle β such that either α is predicated of β and β of γ (first figure), or β is predicated of both α and γ (second figure), or else both α and γ are predicated of β (third figure). All syllogisms must fall into one or another of these figures.

But there is plainly a fourth possibility, that β is predicated of α and γ of β . Many later logicians recognized such syllogisms as belonging to a separate, fourth figure. Aristotle explicitly mentioned such syllogisms but did not group them under a separate figure; his failure to do so has prompted much speculation among commentators and historians. Other logicians included these syllogisms under the first figure. The earliest to do this was Theophrastus (see below), who reinterpreted the first figure in so doing.

Four figures, each with three propositions in one of four forms (A, E, I, O), yield a total of 256 possible syllogistic patterns. Each pattern is called a mood. Only 24 moods are valid, 6 in each figure. Some valid moods may be derived from others by subalternation—that is, if premises validly yield a conclusion of form A, the same premises will yield the corresponding conclusion of form I. So too with forms E and O. Such derived moods were not discussed by Aristotle; they seem to have been first recognized by Ariston of Alexandria (c. 50 bc). In the Middle Ages they were called "subalternate" moods. Disregarding them, there are 4 valid moods in each of the first two figures, 6 in the third figure, and 5 in the fourth. Aristotle recognized all 19 of them.

Here are the valid moods, including subalternate ones, under their medieval mnemonic names (subalternate moods are marked with an asterisk):

- First figure: Barbara, Celarent, Darii, Ferio, *Barbari, *Celaront.
- Second figure: Cesare, Camestres, Festino, Baroco, *Cesaro, *Camestrop.
- Third figure: Darapti, Disamis, Datisi, Felapton, Bocardo, Ferison.
- Fourth figure: Bramantip, Camenes, Dimaris, Fesapo, Fresison, *Camenop.

The sequence of vowels in each name indicates the sequence of categorical propositions in the mood in the order: major, minor, conclusion. Thus, for example, Celarent is a first figure syllogism with an E-form major, A-form minor, and E-form conclusion.

If one assumes the nonsubalternate moods of the first figure, then, with two exceptions, all valid moods in the other figures can be proved by "reducing" them to one of those "axiomatic" first-figure moods. This reduction shows that, if the premises of the reducible mood are true, then it follows, by rules of conversion and one of the axiomatic moods, that the conclusion is true. The procedure is encoded in the medieval names:

1. The initial letter is the initial letter of the first-figure mood to which the given mood is reducible. Thus Felapton is reducible to Ferio.

2. When it is not the final letter, "s" after a vowel means "Convert the sentence simply," and "p" there means "Convert the sentence *per accidens*."
3. When "s" or "p" is the final letter, the conclusion of the first-figure syllogism to which the mood is reduced must be converted simply or *per accidens*, respectively.
4. The letter "m" means "Change the order of the premises."
5. When it is not the first letter, "c" means that the syllogism cannot be directly reduced to the first figure but must be proved by *reductio ad absurdum*. (There are two such moods; see below.)
6. The letters "b" and "d" (except as initial letters) and "l," "n," "t," and "r" serve only to facilitate pronunciation.

Thus the premises of Felapton (third figure) are "No β is an α " and "Every β is a γ ." Convert the minor premise *per accidens* to "Some γ is a β ," as instructed by the "p" after the second vowel. This new proposition and the major premise of Felapton form the premises of a syllogism in Ferio (first figure), the conclusion of which is "Some γ is not an α ," which is also the conclusion of Felapton. Hence, given Ferio and the rule of *per accidens* conversion, the premises of Felapton validly imply its conclusion. In this sense, Felapton has been "reduced" to Ferio.

The two exceptional cases, which must be proven indirectly by *reductio ad absurdum*, are Baroco and Bocardo. Both are reducible indirectly to Barbara in the first figure as follows: Assume the A-form premise (the major in Baroco, the minor in Bocardo). Assume the contradictory of the conclusion. These yield a syllogism in Barbara, the conclusion of which contradicts the O-form premise of the syllogism to be reduced. Thus, given Barbara as axiomatic, and given the premises of the reducible syllogism, the contradictory of its conclusion is false, so that the original conclusion is true.

Reduction and indirect proof together suffice to prove all moods not in the first figure. This fact, which Aristotle himself showed, makes his syllogistic the first deductive system in the history of logic.

While the medieval names of the moods contain a great deal of information, they provide no way by themselves to determine to which figure a mood belongs, and so no way to reconstruct the actual form of the syllogism. Mnemonic verses were developed in the Middle Ages for this purpose.

Categorical propositions in which α is merely said to belong (or not) to some or every β are called assertoric categorical propositions; syllogisms composed solely of such categoricals are called assertoric syllogisms. Aristotle was also interested in categoricals in which α is said to belong (or not) necessarily or possibly to some or every β . Such categoricals are called modal categoricals, and syllogisms in which the component categoricals are modal are called modal syllogisms (they are sometimes called "mixed" if only one of the premises is modal).

Aristotle discussed two notions of the "possible": (1) as what is not impossible (*i.e.*, the opposite of which is not necessary) and (2) as what is neither necessary nor impossible (*i.e.*, the contingent). In his modal syllogistic, the term "possible" (or "contingent") is always used in sense 2 in syllogistic premises, but it is sometimes used in sense 1 in syllogistic conclusions if a conclusion in sense 2 would be incorrect.

Aristotle's procedure in his modal syllogistic is to survey each valid mood of the assertoric syllogistic and then to test the several modal syllogisms that can be formed from an assertoric mood by changing one or more of its component categoricals into a modal categorical. The interpretation of this part of Aristotle's logic, and the correctness of his arguments, have been disputed since antiquity.

Although Aristotle did not develop a full theory of propositions in tenses other than the present, there is a famous passage in the *De interpretatione* that was influential in later developments in this area. In chapter 9 of that work, Aristotle discussed the assertion "There will be a sea battle tomorrow." The discussion assumes that as of now the question is still unsettled. Although there are different interpretations of the passage, Aristotle seems there to have been maintaining that although now, before

The fourth figure

Assertoric and modal syllogisms

Reduction

the fact, it is neither true nor false that there will be a sea battle tomorrow, nevertheless it is true even now, before the fact, that there either will or will not be a sea battle tomorrow. In short, Aristotle appears to have affirmed the law of excluded middle (for any proposition replacing " p ," it is true that either p or not- p), but to have denied the principle of bivalence (that every proposition is either true or false) in the case of future contingent propositions.

Aristotle's logic presupposes several principles that he did not explicitly formulate about logical relations among any propositions whatever, independent of the propositions' internal analyses into categorical or any other form. For example, it presupposes that the principle "If p then q ; but p ; therefore q " (where p and q are replaced by any propositions) is valid. Such patterns of inference belong to what is called the logic of propositions. Aristotle's logic is, by contrast, a logic of terms in the sense described above. A sustained study of the logic of propositions came only after Aristotle.

THEOPHRASTUS OF ERESUS

Aristotle's successor as head of his school at Athens was Theophrastus of Eresus (c. 371–c. 286 BC). All Theophrastus' logical writings are now lost, and much of what was said about his logical views by late ancient authors was attributed to both Theophrastus and his colleague Eudemus, so that it is difficult to isolate their respective contributions.

Theophrastus is reported to have added to the first figure of the syllogism the five moods that others later classified under a fourth figure. These moods were then called indirect moods of the first figure. In order to accommodate them, he had in effect to redefine the first figure as that in which the middle is the subject in one premise and the predicate in the other, not necessarily the subject in the major premise and the predicate in the minor, as Aristotle had it.

Theophrastus' most significant departure from Aristotle's doctrine occurred in modal syllogistic. He abandoned Aristotle's notion of the possible as neither necessary nor impossible and adopted Aristotle's alternative notion of the possible as simply what is not impossible. This allowed him to effect a considerable simplification in Aristotle's modal theory. Thus, his conversion laws for modal categoricals were exact parallels to the corresponding laws for assertoric categoricals. In particular, for Theophrastus "problematic" universal negatives ("No β is possibly an α ") can be simply converted. Aristotle had denied this.

In addition, Theophrastus adopted a rule that the conclusion of a valid modal syllogism can be no stronger than its weakest premise. (Necessity is stronger than possibility, and an assertoric claim without any modal qualification is intermediate between the two). This rule simplifies modal syllogistic and eliminates several moods that Aristotle had accepted. Yet Theophrastus himself allowed certain modal moods that, combined with the principle of indirect proof (which he likewise accepted), yield results that perhaps violate this rule.

Theophrastus also developed a theory of inferences involving premises of the form " α is universally predicated of everything of which γ is universally predicated" and of related forms. Such propositions he called proleptic propositions, and inferences involving them were termed proleptic syllogisms. Greek *prolepsis* can mean "something taken in addition," and Theophrastus claimed that propositions like these implicitly contain a third, indefinite term, in addition to the two definite terms (" α " and " γ " in the example).

The term proleptic proposition appears to have originated with Theophrastus, although Aristotle discussed such propositions briefly in his *Prior Analytics* without exploring their logic in detail. The implicit third term in a proleptic proposition Theophrastus called the middle. After an analogy with syllogistic for categorical propositions, he distinguished three "figures" for proleptic propositions and syllogisms, based on the position of the implicit middle. The proleptic proposition " α is universally predicated of everything that is universally predicated of γ " belongs to the first figure and can be a premise in a first-figure

proleptic syllogism. "Everything predicated universally of α is predicated universally of γ " belongs to the second figure and can be a premise in a second-figure syllogism, and so too " α is universally predicated of everything of which γ is universally predicated" for the third figure. Thus, for example, the following is a proleptic syllogism in the third figure: " α is universally affirmed of everything of which γ is universally affirmed; γ is universally affirmed of β ; therefore, α is universally affirmed of β ."

Theophrastus observed that certain proleptic propositions are equivalent to categoricals and differ from them only "potentially" or "verbally." Some late ancient authors claimed that this made proleptic syllogisms superfluous. But in fact not all proleptic propositions are equivalent to categoricals.

Theophrastus is also credited with investigations into hypothetical syllogisms. A hypothetical proposition, for Theophrastus, is a proposition made up of two or more component propositions (e.g., " p or q ," or "if p then q "), and a hypothetical syllogism is an inference containing at least one hypothetical proposition as a premise. The extent of Theophrastus' work in this area is uncertain, but it appears that he investigated a class of inferences called totally hypothetical syllogisms, in which both premises and the conclusion are conditionals. This class would include, for example, syllogisms such as "If α then β ; if β then γ ; therefore, if α then γ ," or "if α then β ; if not α then γ ; therefore, if not β then γ ." As with his proleptic syllogisms, Theophrastus divided these totally hypothetical syllogisms into three "figures," after an analogy with categorical syllogistic.

Theophrastus was the first person in the history of logic known to have examined the logic of propositions seriously. Still, there was no sustained investigation in this area until the period of the Stoics.

THE MEGARIANS AND STOICS

Throughout the ancient world, the logic of Aristotle and his followers was one main stream. But there was also a second tradition of logic, that of the Megarians and Stoics.

The Megarians were followers of Euclid (or Euclides) of Megara (c. 430–c. 360 BC), a pupil of Socrates. In logic the most important Megarians were Diodorus Cronus (4th century BC) and his pupil Philo of Megara. The Stoics were followers of Zeno of Citium (c. 336–c. 265 BC). By far the most important Stoic logician was Chrysippus (c. 279–206 BC). The influence of Megarian on Stoic logic is indisputable, but many details are uncertain, since all but fragments of the writings of both groups are lost.

The Megarians were interested in logical puzzles. Many paradoxes have been attributed to them, including the "liar paradox" (someone says that he is lying; is his statement true or false?), the discovery of which has sometimes been credited to Eubulides of Miletus, a pupil of Euclid of Megara. The Megarians also discussed how to define various modal notions and debated the interpretation of conditional propositions.

Diodorus Cronus originated a mysterious argument called the Master Argument. It claimed that the following three propositions are jointly inconsistent, so that at least one of them is false:

1. Everything true about the past is now necessary. (That is, the past is now settled, and there is nothing to be done about it.)
2. The impossible does not follow from the possible.
3. There is something that is possible, and yet neither is nor will be true. (That is, there are possibilities that will never be realized.)

It is unclear exactly what inconsistency Diodorus saw among these propositions. Whatever it was, Diodorus was unwilling to give up 1 or 2, and so rejected 3. That is, he accepted the opposite of 3, namely: Whatever is possible either is or will be true. In short, there are no possibilities that are not realized now or in the future. It has been suggested that the Master Argument was directed against Aristotle's discussion of the sea battle tomorrow in the *De interpretatione*.

Diodorus also proposed an interpretation of conditional propositions. He held that the proposition "If p , then q "

Hypothetical syllogisms

Contributions to modal syllogistic

The Master Argument

is true if and only if it neither is nor ever was possible for the antecedent p to be true and the consequent q to be false simultaneously. Given Diodorus' notion of possibility, this means that a true conditional is one that at no time (past, present, or future) has a true antecedent and a false consequent. Thus, for Diodorus a conditional does not change its truth value; if it is ever true, it is always true. But Philo of Megara had a different interpretation. For him, a conditional is true if and only if it does not now have a true antecedent and a false consequent. This is exactly the modern notion of material implication. In Philo's view, unlike Diodorus', conditionals may change their truth value over time.

These and other theories of modality and conditionals were discussed not only by the Megarians but by the Stoics as well. Stoic logicians, like the Megarians, were not especially interested in scientific demonstration in Aristotle's special sense. They were more concerned with logical issues arising from debate and disputation: fallacies, paradoxes, forms of refutation. Aristotle had also written about such things, but his interests gradually shifted to his special notion of science. The Stoics kept their interest focused on disputation and developed their studies in this area to a high degree.

Unlike the Aristotelians, the Stoics developed propositional logic to the neglect of term logic. They did not produce a system of logical laws arising from the internal structure of simple propositions, as Aristotle had done with his account of opposition, conversion, and syllogistic for categorical propositions. Instead, they concentrated on inferences from hypothetical propositions as premises. Theophrastus had already taken some steps in this area, but his work had little influence on the Stoics.

Stoic logicians studied the logical properties and defining features of words used to combine simpler propositions into more complex ones. In addition to the conditional, which had already been explored by the Megarians, they investigated disjunction ("or") and conjunction ("and"), along with words like "since" and "because." Some of these they defined truth-functionally (*i.e.*, solely in terms of the truth or falsehood of the propositions they combined). For example, they defined a disjunction as true if and only if exactly one disjunct is true (the modern "exclusive" disjunction). They also knew "inclusive" disjunction (defined as true when at least one disjunct is true), but this was not widely used. More important, the Stoics seem to have been the first to show how some of these truth-functional words may be defined in terms of others.

Unlike Aristotle, who typically formulated his syllogisms as conditional propositions, the Stoics regularly presented principles of logical inference in the form of schematic arguments. While Aristotle had used Greek letters as variables replacing terms, the Stoics used ordinal numerals as variables replacing whole propositions. Thus: "Either the first or the second; but not the second; therefore, the first." Here the expressions "the first" and "the second" are variables or placeholders for propositions, not terms.

Chrysippus regarded five valid inference schemata as basic or indemonstrable. They are:

1. If the first, then the second; but the first; therefore, the second.
2. If the first, then the second; but not the second; therefore not the first.
3. Not both the first and the second; but the first; therefore, not the second.
4. Either the first or the second; but the first; therefore, not the second.
5. Either the first or the second; but not the second; therefore, the first.

Using these five "indemonstrables," Chrysippus proved the validity of many further inference schemata. Indeed, the Stoics claimed (falsely, it seems) that all valid inference schemata could be derived from the five indemonstrables.

The differences between Aristotelian and Stoic logic were ones of emphasis, not substantive theoretical disagreements. At the time, however, it appeared otherwise. Perhaps because of their real disputes in other areas, Aristotelians and Stoics at first saw themselves as holding incompatible theories in logic as well. But by the

late 1st century BC, an eclectic movement had begun to weaken these hostilities. Thereafter the two traditions were combined in commentaries and handbooks for general education.

LATE REPRESENTATIVES OF ANCIENT GREEK LOGIC

After Chrysippus, little important logical work was done in Greek. But the commentaries and handbooks that were written did serve to consolidate the previous traditions and in some cases are the only extant sources for the doctrines of earlier writers. Among late authors, Galen the physician (AD 129–c. 199) wrote several commentaries, now lost, and an extant *Introduction to Dialectic*. Galen observed that the study of mathematics and logic was important to a medical education, a view that had considerable influence in the later history of logic, particularly in the Arab world. Tradition has credited Galen with "discovering" the fourth figure of the Aristotelian syllogism, although in fact he explicitly rejected it.

Alexander of Aphrodisias (fl. c. AD 200) wrote extremely important commentaries on Aristotle's writings, including the logical works. Other important commentators include Porphyry of Tyre (c. 232–before 306), Ammonius Hermioui (5th century), Simplicius (6th century), and John Philoponus (6th century). Sextus Empiricus (late 2nd–early 3rd centuries) and Diogenes Laërtius (probably early 3rd century) are also important sources for earlier writers. Significant contributions to logic were not made again in Europe until the 12th century.

Medieval logic

TRANSMISSION OF GREEK LOGIC TO THE LATIN WEST

As the Greco-Roman world disintegrated and gave way to the Middle Ages, knowledge of Greek declined in the West. Nevertheless, several authors served as transmitters of Greek learning to the Latin world. Among the earliest of them, Cicero (106–43 BC) introduced Latin translations for technical Greek terms. Although his translations were not always finally adopted by later authors, he did make it possible to discuss logic in a language that had not previously had any precise vocabulary for it. In addition, he preserved much information about the Stoics. In the 2nd century AD Lucius Apuleius passed on some knowledge of Greek logic in his *De philosophia rationali* ("On Rational Philosophy").

In the 4th century Marius Victorinus produced Latin translations of Aristotle's *Categories* and *De interpretatione* and of Porphyry of Tyre's *Isagoge* ("Introduction," on Aristotle's *Categories*), although these translations were not very influential. He also wrote logical treatises of his own. A short *De dialectica* ("On Dialectic"), doubtfully attributed to St. Augustine (354–430), shows evidence of Stoic influence, although it had little influence of its own. The pseudo-Augustinian *Decem categoriae* ("Ten Categories") is a late 4th-century Latin paraphrase of a Greek compendium of the *Categories*. In the late 5th century Martianus Capella's allegorical *De nuptiis Philologiae et Mercurii* (*The Marriage of Philology and Mercury*) contains "On the Art of Dialectic" as book IV.

The first truly important figure in medieval logic was Boethius (480–524/525). Like Victorinus, he translated Aristotle's *Categories* and *De interpretatione* and Porphyry's *Isagoge*, but his translations were much more influential. He also seems to have translated the rest of Aristotle's *Organon*, except for the *Posterior Analytics*, but the history of those translations and their circulation in Europe is much more complicated; they did not come into widespread use until the first half of the 12th century. In addition, Boethius wrote commentaries and other logical works that were of tremendous importance throughout the Latin Middle Ages. Until the 12th century his writings and translations were the main sources for medieval Europe's knowledge of logic. In the 12th century they were known collectively as the *Logica vetus* ("Old Logic").

ARABIC LOGIC

Between the time of the Stoics and the revival of logic in 12th-century Europe, the most important logical work

Aristotle's
commenta-
tors

Boethius

Disjunc-
tion and
conjunc-
tion

was done in the Arab world. Arabic interest in logic lasted from the 9th to the 16th century, although the most important writings were done well before 1300.

Syrian Christian authors in the late 8th century were among the first to introduce Alexandrian scholarship to the Arab world. Through Galen's influence, these authors regarded logic as important to the study of medicine. (This link with medicine continued throughout the history of Arabic logic and, to some extent, later in medieval Europe.) By about 850, at least Porphyry's *Isagoge* and Aristotle's *Categories*, *De interpretatione*, and *Prior Analytics* had been translated via Syriac into Arabic. Between 830 and 870 the philosopher and scientist al-Kindi (c. 805–873) produced in Baghdad what seem to have been the first Arabic writings on logic that were not translations. But these writings, now lost, were probably mere summaries of others' work.

By the late 9th century, the school of Baghdad was the focus of logic studies in the Arab world. Most of the members of this school were Nestorian or Jacobite Christians, but the Muslim al-Fārābī (c. 873–950) wrote important commentaries and other logical works there that influenced all later Arabic logicians. Many of these writings are now lost, but among the topics al-Fārābī discussed were future contingents (in the context of Aristotle's *De interpretatione*, chapter 9), the number and relation of the categories, the relation between logic and grammar, and non-Aristotelian forms of inference. This last topic showed the influence of the Stoics. Al-Fārābī, along with Avicenna and Averroës (see below), was among the best logicians the Arab world produced.

By 1050 the school of Baghdad had declined. The 11th century saw very few Arabic logicians, with one distinguished exception: the Persian Ibn Sinā, or Avicenna (980–1037), perhaps the most original and important of all Arabic logicians. Avicenna abandoned the practice of writing on logic in commentaries on the works of Aristotle and instead produced independent treatises. He sharply criticized the school of Baghdad for what he regarded as their slavish devotion to Aristotle. Among the topics Avicenna investigated were quantification of the predicates of categorical propositions, the theory of definition and classification, and an original theory of "temporally modalized" syllogistic, in which premises include such modifiers as "at all times," "at most times," and "at some time."

The Persian mystic and theologian al-Ghazālī, or Algazel, (1058–1111), followed Avicenna's logic, although he differed sharply from Avicenna in other areas. Al-Ghazālī was not a significant logician but is important nonetheless because of his influential defense of the use of logic in theology.

In the 12th century the most important Arab logician was Ibn Rushd, or Averroës (1126–98). Unlike the Persian followers of Avicenna, Averroës worked in Moorish Spain, where he revived the tradition of al-Fārābī and the school of Baghdad by writing penetrating commentaries on Aristotle's works, including the logical ones. Such was the stature of these excellent commentaries that, when they were translated into Latin in the 1220s or 1230s, Averroës was often referred to simply as "the Commentator."

After Averroës, logic declined in western Islām because of the antagonism felt to exist between logic and philosophy on the one hand and Muslim orthodoxy on the other. But in eastern Islām, owing in part to the work of al-Ghazālī, logic was not regarded as being so closely linked with philosophy. Instead, it was viewed as a tool that could be profitably used in any field of study, even (as al-Ghazālī had done) on behalf of theology against the philosophers. Thus the logical tradition continued in Persia long after it died out in Spain. The 13th century produced a large number of logical writings, but these were mostly unoriginal textbooks and handbooks. After about 1300, logical study was reduced to producing commentaries on these earlier, already derivative handbooks.

THE REVIVAL OF LOGIC IN EUROPE

St. Anselm and Peter Abelard. Except in the Arabic world, there was little activity in logic between the time of Boethius and the 12th century. Certainly Byzantium

produced nothing of note. In Latin Europe there were a few authors, including Alcuin of York (c. 730–804) and Garland the Computist (*fl.* c. 1040). But it was not until late in the 11th century that serious interest in logic revived. St. Anselm of Canterbury (1033–1109) discussed semantical questions in his *De grammatico*, and investigated the notions of possibility and necessity in surviving fragments, but these texts did not have much influence. More important was Anselm's general method of using logical techniques in theology. His example set the tone for much that was to follow.

The first important Latin logician after Boethius was Peter Abelard (1079–1142). He wrote three sets of commentaries and glosses on Porphyry's *Isagoge* and Aristotle's *Categories* and *De interpretatione*; these were the *Introductiones parvulorum* (also containing glosses on some writings of Boethius), *Logica "Ingredientibus,"* and *Logica "Nostrorum petitioni sociorum"* (on the *Isagoge* only), together with the independent treatise *Dialectica* (extant in part). These works show a familiarity with Boethius but go far beyond him. Among the topics discussed insightfully by Abelard are the role of the copula in categorical propositions, the effects of different positions of the negation sign in categorical propositions, modal notions like "possibility," future contingents (as treated, for example, in chapter 9 of Aristotle's *De interpretatione*), and conditional propositions or "consequences."

Abelard's fertile investigations raised logical study in medieval Europe to a new level. His achievement is all the more remarkable since the sources at his disposal were the same ones that had been available in Europe for the preceding 600 years: Aristotle's *Categories* and *De interpretatione* and Porphyry's *Isagoge*, together with the commentaries and independent treatises by Boethius.

The "properties of terms" and discussions of fallacies. Even in Abelard's lifetime, however, things were changing. After about 1120, Boethius' translations of Aristotle's *Prior Analytics*, *Topics*, and *Sophistic Refutations* began to circulate. Sometime in the second quarter of the 12th century, James of Venice translated the *Posterior Analytics* from Greek, thus making the whole of the *Organon* available in Latin. These newly available Aristotelian works were known collectively as the *Logica nova* ("New Logic"). In a flurry of activity, others in the 12th and 13th centuries produced additional translations of these works and of Greek and Arabic commentaries on them, along with many other philosophical writings and other works from Greek and Arabic sources.

The *Sophistic Refutations* proved an important catalyst in the development of medieval logic. It is a little catalog of fallacies, how to avoid them, and how to trap others into committing them. The work is very sketchy. Many kinds of fallacies are not discussed, and those that are could have been treated differently. Unlike the *Posterior Analytics*, the *Sophistic Refutations* was relatively easy to understand. And unlike the *Prior Analytics*—where, except for modal syllogistic, Aristotle had left little to be done—there was obviously still much to be investigated about fallacies. Moreover, the discovery of fallacies was especially important in theology, particularly in the doctrines of the Trinity and the Incarnation. In short, the *Sophistic Refutations* was tailor-made to exercise the logical ingenuity of the 12th century. And that is exactly what happened.

The *Sophistic Refutations*, and the study of fallacy it generated, produced an entirely new logical literature. A genre of *sophismata* ("sophistical") treatises developed that investigated fallacies in theology, physics, and logic. The theory of "supposition" (see below) also developed out of the study of fallacies. Whole new kinds of treatises were written on what were called "the properties of terms," semantic properties important in the study of fallacy. In addition, a new genre of logical writings developed on the topic of "syncategoremata"—expressions such as "only," "inasmuch as," "besides," "except," "lest," and so on, which posed quite different logical problems than did the terms and logical particles in traditional categorical propositions or in the simpler kind of "hypothetical" propositions inherited from the Stoics. The study of valid

Peter
Abelard

Avicenna

Sophismata
treatises

inference generated a literature on "consequences" that went into far more detail than any previous studies. By the late 12th or early 13th century, special treatises were devoted to *insolubilia* (semantic paradoxes such as the liar paradox, "This sentence is false") and to a kind of disputation called "obligationes," the exact purpose of which is still in question.

All these treatises, and the logic contained in them, constitute the peculiarly medieval contribution to logic. It is primarily on these topics that medieval logicians exercised their best ingenuity. Such treatises, and their logic, were called the *Logica moderna* ("Modern Logic"), or "terminist" logic, because they laid so much emphasis on the "properties of terms." These developments began in the mid-12th century, and continued to the end of the Middle Ages.

DEVELOPMENTS IN THE 13TH AND EARLY 14TH CENTURIES

In the 13th century the *sophismata* literature continued and deepened. In addition several authors produced summary works that surveyed the whole field of logic, including the "Old" and "New" logic as well as the new developments in the *Logica moderna*. These compendia are often called "*summulae*" ("little summaries"), and their authors "summulists." Among the most important of the summulists are: (1) Peter of Spain (also known as Petrus Hispanus; later Pope John XXI), who wrote a *Tractatus* more commonly known as *Summulae logicales* ("Little Summaries of Logic") probably in the early 1230s; it was used as a textbook in some late medieval universities; (2) Lambert of Auxerre, who wrote a *Logica* sometime between 1253 and 1257; and (3) William of Sherwood, who produced *Introductiones in logicam* (*Introduction to Logic*) and other logical works sometime about the mid-century.

Despite his significance in other fields, Thomas Aquinas is of little importance in the history of logic. He did write a treatise on modal propositions and another one on fallacies. But there is nothing especially original in these works; they are early writings and are confined to passing on received doctrine. He also wrote an incomplete commentary on the *De interpretatione*, but it is of no great logical significance.

About the end of the 13th century John Duns Scotus (c. 1266–1308) composed several works on logic. There also are some very interesting logical texts from the same period that have been falsely attributed to Scotus and were published in the 17th century among his authentic works. These are now referred to as the works of "the Pseudo-Scotus," although they may not all be by the same author.

The first half of the 14th century saw the high point of medieval logic. Much of the best work was done by people associated with the University of Oxford. Among them were William of Ockham (c. 1285–1347), the author of an important *Summa logicae* ("Summary of Logic") and other logical writings. Perhaps because of his importance in other areas of medieval thought, Ockham's originality in logic has sometimes been exaggerated. But there is no doubt that he was one of the most important logicians of the century. Another Oxford logician was Walter Burley (or Burleigh), an older contemporary of Ockham. Burley was a bitter opponent of Ockham in metaphysics. He wrote a work *De puritate artis logicae* ("On the Purity of the Art of Logic"; in two versions), apparently in response and opposition to Ockham's views, although on some points Ockham simply copied Burley almost verbatim.

Slightly later, on the Continent, Jean Buridan was a very important logician at the University of Paris. He wrote mainly during the 1330s and '40s. In many areas of logic and philosophy his views were close to Ockham's, although the extent of Ockham's influence on Buridan is not clear. Buridan's *Summulae de dialectica* ("Little Summaries of Dialectic"), intended for instructional use at Paris, was largely an adaptation of Peter of Spain's *Summulae logicales*. He appears to have been the first to use Peter of Spain's text in this way. Originally meant as the last treatise of his *Summulae de dialectica*, Buridan's extremely interesting *Sophismata* (published separately in early editions) discusses many issues in semantics and philosophy of logic. Among Buridan's pupils was Albert

of Saxony (d. 1390), the author of a *Perutilis logica* ("A Very Useful Logic") and later first rector of the University of Vienna. Albert was not an especially original logician, although his influence was by no means negligible.

The theory of supposition. Many of the characteristically medieval logical doctrines in the *Logica moderna* centred around the notion of "supposition" (*suppositio*). Already by the late 12th century the theory of supposition had begun to form. In the 13th century, special treatises on the topic multiplied. The summulists all discussed it at length. Then, after about 1270, relatively little was heard about it. In France, supposition theory was replaced by a theory of "speculative grammar" or "modism" (so called because it appealed to "modes of signifying"). Modism was not so popular in England, but there too the theory of supposition was largely neglected in the late 13th century. In the early 14th century the theory reemerged both in England and on the Continent. Burley wrote a treatise on the topic in about 1302, and Buridan revived the theory in France in the 1320s. Thereafter the theory remained the main vehicle for semantic analysis until the end of the Middle Ages.

Supposition theory, at least in its 14th-century form, is best viewed as two theories under one name. The first, sometimes called the theory of "supposition proper," is a theory of reference and answers the question "To what does a given occurrence of a term refer in a given proposition?" In general (the details depend on the author) three main types of supposition were distinguished: (1) personal supposition (which, despite the name, need not have anything to do with persons), (2) simple supposition, and (3) material supposition. These types are illustrated, respectively, by the occurrences of the term *horse* in the statements "Every horse is an animal" (in which the term *horse* refers to individual horses), "Horse is a species" (in which the term refers to a universal), and "Horse is a monosyllable" (in which it refers to the spoken or written word). The theory was elaborated and refined by considering how reference may be broadened by tense and modal factors (for example, the term *horse* in "Every horse will die," which may refer to future as well as present horses) or narrowed by adjectives or other factors (for example, *horse* in "Every horse in the race is less than two years old").

The second part of supposition theory applies only to terms in personal supposition. It divides personal supposition into several types, including (again the details vary according to the author): (1) determinate (e.g., *horse* in "Some horse is running"), (2) confused and distributive (e.g., *horse* in "Every horse is an animal"), and (3) merely confused (e.g., *animal* in "Every horse is an animal"). These types were described in terms of a notion of "descent to (or ascent from) singulars." For example, in the statement "Every horse is an animal," one can "descend" under the term "horse" to: "This horse is an animal, and that horse is an animal, and so on," but one cannot validly "ascend" from "This horse is an animal" to the original proposition. There are many refinements and complications.

The purpose of this second part of the theory of supposition has been disputed. Since the question of what it is to which a given occurrence of a term refers is already answered in the first part of supposition theory, the purpose of this second part must have been different. The main suggestions are (1) that it was devised to help detect and diagnose fallacies, (2) that it was intended as a theory of truth conditions for propositions or as a theory of analyzing the senses of propositions, and (3) that, like the first half of supposition theory, it originated as part of an account of reference, but, once its theoretical insufficiency for that task was recognized, it was gradually divorced from that first part of supposition theory and by the early 14th century was left as a conservative vestige that continued to be disputed but no longer had any question of its own to answer. There are difficulties with all of these suggestions. The theory of supposition survived beyond the Middle Ages and was frequently applied not only in logical discussions but also in theology and in the natural sciences.

The theory of "supposition proper"

Connotation

In addition to supposition and its satellite theories, several logicians during the 14th century developed a sophisticated theory of "connotation" (*connotatio* or *appellatio*; in which the term *black*, for instance, not only refers to black things but also "connotes" the quality, blackness, that they possess) and a subtle theory of "mental language," in which tools of semantic analysis were applied to epistemology and the philosophy of mind. Important treatises on *insolubilia* and *obligationes*, as well as on the theory of consequence or inference, continued to be produced in the 14th century, although the main developments there were completed by mid-century.

Developments in modal logic. Medieval logicians continued the tradition of modal syllogistic inherited from Aristotle. In addition, modal factors were incorporated into the theory of supposition. But the most important developments in modal logic occurred in three other contexts: (1) whether propositions about future contingent events are now true or false (Aristotle had raised this question in *De interpretatione*, chapter 9), (2) whether a future contingent event can be known in advance, and (3) whether God (who, the tradition says, cannot be acted upon causally) can know future contingent events. All these issues link logical modality with time. Thus, Peter Aureoli (c. 1280–1322) held that if something is in fact φ (' φ ' is some predicate) but can be not- φ , then it is capable of changing from being φ to being not- φ .

Duns Scotus in the late 13th century was the first to sever the link between time and modality. He proposed a notion of possibility that was not linked with time but based purely on the notion of semantic consistency. This radically new conception had a tremendous influence on later generations down to the 20th century. Shortly afterward, Ockham developed an influential theory of modality and time that reconciles the claim that every proposition is either true or false with the claim that certain propositions about the future are genuinely contingent.

LATE MEDIEVAL LOGIC

Most of the main developments in medieval logic were in place by the mid-14th century. On the Continent, the disciples of Jean Buridan—Albert of Saxony (c. 1316–90), Marsilius of Inghen (d. 1399), and others—continued and developed the work of their predecessors. In 1372 Pierre d'Ailly wrote an important work, *Conceptus et insolubilia* (*Concepts and Insolubles*), which appealed to a sophisticated theory of mental language in order to solve semantic paradoxes such as the liar paradox.

In England the second half of the 14th century produced several logicians who consolidated and elaborated earlier developments. Their work was not very original, although it was often extremely subtle. Many authors during this period compiled brief summaries of logical topics intended as textbooks. The doctrine in these little summaries is remarkably uniform, making it difficult to determine who their authors were. By the early 15th century, informal collections of these treatises had been gathered under the title *Libelli sophistarum* ("Little Books for Arguers")—one collection for Oxford and a second for Cambridge; both were printed in early editions. Among the notable logicians of this period are Henry Hopton (fl. 1357), John Wycliffe (c. 1330–84), Richard Lavenham (d. after 1399), Ralph Strode (fl. c. 1360), Richard Ferrybridge (or Feribrigge; fl. c. 1360s), and John Venator (also known as John Huntman or Hunter; fl. 1373).

Beginning in 1390, the Italian Paul of Venice studied for at least three years at Oxford and then returned to teach at Padua and elsewhere in Italy. Although English logic was studied in Italy even before Paul's return, his own writings advanced this study greatly. Among Paul's logical works were the very popular *Logica parva* ("Little Logic"), printed in several early editions, and possibly the huge *Logica magna* ("Big Logic") that has sometimes been regarded as a kind of encyclopaedia of the whole of medieval logic.

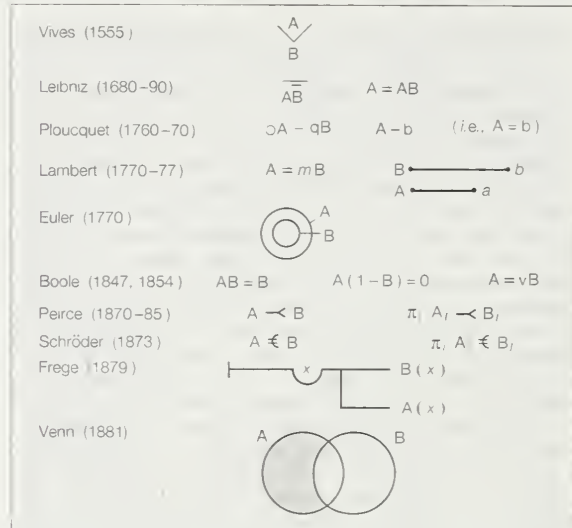
After about 1400, serious logical study was dead in England. However, it continued to be pursued on the Continent until the end of the Middle Ages and afterward.

(P.V.S.)

Libelli
sophistarum

Modern logic

It is customary to speak of logic since the Renaissance as "modern logic." This is not to suggest that there was a smooth development of a unified conception of reasoning, or that the logic of this period is "modern" in the usual sense. Logic in the modern era has exhibited an extreme diversity, and its chaotic development has reflected all too clearly the surrounding political and intellectual turmoil. These upheavals include the Renaissance itself, the diminishing role of the Roman Catholic church and of Latin, the Reformation and subsequent religious wars, the scientific revolution and the growth of modern mathematics, the rise and fall of empires and nation-states, and the waxing influence of the New World and the former Soviet Union.



Representations of the universal affirmative, "All A's are B's" in modern logic.

THE 16TH CENTURY

Renaissance writers sometimes denounced all of scholastic logic. The humanism of the Renaissance is often seen as promoting the study of Greek and Roman classics, but Aristotle's logic was frequently regarded as being so hopelessly bound together with "sterile" medieval logic as to constitute an exception to this spirit of rebirth. Some, such as Martin Luther (1483–1546), were repelled by any hint of Aristotelianism. Others, such as the great humanist essayist Desiderius Erasmus (1466–1536), occasionally praised Aristotle but never his logical theory; like many writers in the Renaissance, Erasmus found in the theory of the syllogism only "subtlety and arid ingenuity" (Johan Huizinga, *Erasmus* [1924]). The German Lutheran humanist Philipp Melanchthon (1497–1560) had a more balanced appreciation of Aristotle's logic. Melanchthon's *Compendaria dialectices ratio* ("Brief Outline of Dialectics") of 1520, built upon his *Institutiones Rhetoricae* of the previous year, became a popular Lutheran text. There he described his purpose as presenting "a true, pure and uncomplicated logic, just as we have received it from Aristotle and some of his judicious commentators." Elsewhere, influential writers such as Rabalais, Petrarch, and Montaigne had few kind words for logic as they knew it.

The French reformer and pamphleteer Petrus Ramus (Pierre de la Ramée) was also the author of extremely influential "Reform" logical texts. His *Dialectique* (Dialectics) of 1555 (translated into English in 1574) was the first major logical work in a modern language. In this work and in his *Dialecticae libri duo* ("Two Books of Dialectics") of 1556 he combined attacks on scholastic logic, an emphasis on the use of logic in actual arguments ("dialectics"), and a presentation of a much simplified approach to categorical syllogism (without an attempt to follow Aristotle). Elsewhere, he proposed that reasoning should be taught by using Euclid's *Elements* rather than by the study of the syllogism. He devoted special attention to valid syllogisms with singular premises, such as "Octavius is the

Petrus
Ramus

heir of Caesar. I am Octavius. Therefore, I am the heir of Caesar." Singular terms (such as proper names) had been treated by earlier logicians: Pseudo-Scotus, among others, had proposed assimilating them to universal propositions by understanding "Julius Caesar is mortal" as "All Julius Caesars are mortal." Although Ramus' proposals for singular terms were not widely accepted, his concern for explicitly addressing them and his refusal to use artificial techniques to convert them to standard forms prefigured more recent interests. Although it had its precursors in medieval semantic thought, Ramus' division of thought into a hierarchy composed of concepts, judgments, arguments, and method was influential in the 17th and 18th centuries.

Scholastic logic remained alive, especially in predominantly Roman Catholic universities and countries, such as Italy and Spain. Some of this work had considerable value, even though it was outside of the mainstream logical tradition, from which it diverged in the 16th century. If the Reform tradition of Melanchthon and Ramus represents one major tradition in modern logic, and the neo-scholastic tradition another, then (here following the historian of logic Nikolai Ivanovich Styazhkin) a third tradition is found in the followers of the Spanish (Majorcan) soldier, priest, missionary, and mystic Ramón Lull (1235–1315). His *Ars magna, generalis et ultima* (1501; "Great, General and Ultimate Art") represents an attempt to symbolize concepts and derive propositions that form various combinations of possibilities. These notions, associated with lore of the Kabbala, later influenced Pascal and Leibniz and the rise of probability theory. Lull's influence can be seen more directly in the work of his fellow Spaniard Juan Luis Vives (1492–1540), who used a V-shaped symbol to indicate the inclusion of one term in another. Other work inspired by Lull includes the logic and notational system of the German logician Johann Heinrich Alsted (1588–1638). The work of Vives and Alsted represents perhaps the first systematic effort at a logical symbolism.

Interest in the systematic symbolization of logic

With the 17th century came increasing interest in symbolizing logic. These symbolizations sometimes took graphic or pictorial forms but more often used letters in the manner of algebra to stand for propositions, concepts, classes, properties, and relations, as well as special symbols for logical notions. Inspired by the triumphs achieved in mathematics after it had turned to the systematic use of special symbols, logicians hoped to imitate this success. The systematic application of symbols and abbreviations and the conscious hope that through this application great progress could be made have been a distinguishing characteristic of modern logic into the 20th century.

The modern era saw major changes not only in the external appearance of logical writings but also in the purposes of logic. Logic for Aristotle was a theory of ideal human reasoning and inference that also had clear pedagogical value. Early modern logicians stressed what they called "dialectics" (or "rhetoric"), because "logic" had come to mean an elaborate scholastic theory of reasoning that was not always directed toward improving reasoning. A related goal was to extend the scope of human reasoning beyond textbook syllogistic theory and to acknowledge that there were important kinds of valid inference that could not be formulated in traditional Aristotelian syllogistic. But another part of the rejection of Aristotelian logic (broadly conceived to include scholastic logic) is best explained by the changing and quite new goals that logic took on in the modern era. One such goal was the development of an ideal logical language that naturally expressed ideal thought and was more precise than natural languages. Another goal was to develop methods of thinking and discovery that would accelerate or improve human thought or would allow its replacement by mechanical devices. Whereas Aristotelian logic had seen itself as a tool for training "natural" abilities at reasoning, later logics proposed vastly improving meagre and wavering human tendencies and abilities. The linking of logic with mathematics was an especially characteristic theme in the modern era. Finally, in the modern era came an intense consciousness of the importance of logical form (forms of sentences, as well as forms or patterns of arguments). Although the medievals made many distinctions among

patterns of sentences and arguments, the modern logical notion of "form" perhaps first crystallized in the work of Sir William Rowan Hamilton and the English mathematician and logician Augustus De Morgan (*De Morgan's Formal Logic* of 1847). The now standard discussions of validity, invalidity, and the self-conscious separation of "formal" from nonformal aspects of sentences and arguments all trace their roots to this work.

THE 17TH CENTURY

The *Logica Hamburgensis* (1638) of Joachim Jung (also called Jungius or Junge) was one replacement for the "Protestant" logic of Melanchthon. Its chief virtue was the care with which late medieval theories and techniques were gathered and presented. Jung devoted considerable attention to valid arguments that do not fit into simpler, standard conceptions of the syllogism and immediate inference. Of special interest is his treatment of quantified relational arguments, then called "oblique" syllogisms because of the oblique (non-nominative) case that is used to express them in Latin. An example is: "The square of an even number is even; 6 is even; therefore, the square of 6 is even." The technique of dealing with such inferences involved rewriting a premise so that the term in the oblique case (for example, "of an even number") would occur in the subject position and thus be amenable to standard syllogistic manipulation. Such arguments had in fact been noticed by Aristotle and were also treated in late medieval logic.

"Oblique" syllogisms

An especially widely used text of the 17th century is usually termed simply the *Port-Royal Logic* after the seat of the anticlerical Jansenist movement outside Paris. It was written by Antoine Arnauld and Pierre Nicole, possibly with others, and was published in French in 1662 with the title *La Logique ou l'art de penser* "Logic or the Art of Thinking". It was promptly translated into Latin and English and underwent many reprintings in the late 17th and 18th centuries. In its outline, it followed Ramus' outline of concept, judgment, argument, and method; it also briefly mentioned oblique syllogisms. The *Port-Royal Logic* followed the general Reform program of simplifying syllogistic theory, reducing the number of syllogistic figures from four, and minimizing distinctions thought to be useless. In addition, the work contained an important contribution to semantics in the form of the distinction between comprehension and extension. Although medieval semantic theory had used similar notions, the Port-Royal notions found their way into numerous 18th- and 19th-century discussions of the meanings and reference of terms; they appeared, for example, in John Stuart Mill's influential text *A System of Logic* (1843). The "comprehension" of a term consisted of all the essential attributes in it (those that cannot be removed without "destroying" the concept), and the extension consisted of all those objects to which the concept applies. Thus the comprehension of the term "triangle" might include the attributes of being a polygon, three-sided, three-angled, and so on. Its extension would include all kinds of triangles. The *Port-Royal Logic* also contained an influential discussion of definitions that was inspired by the work of the French mathematician and philosopher Blaise Pascal. According to this discussion, some terms could not be defined ("primitive" terms), and definitions were divided between nominal and real ones. Real definitions were descriptive and stated the essential properties in a concept, while nominal definitions were creative and stipulated the conventions by which a linguistic term was to be used.

Nominal and real definitions

Discussions of "nominal" and "real" definitions go back at least to the nominalist/realist debates of the 14th century; Pascal's application of the distinction is interesting for the emphasis that it laid on mathematical definitions being nominal and on the usefulness of nominal definitions. Although the Port-Royal logic itself contained no symbolism, the philosophical foundation for using symbols by nominal definitions was nevertheless laid.

One intriguing 17th-century treatment of logic in terms of demonstrations, postulates, and definitions in a Euclidean fashion occurs in the otherwise quite traditional *Logica Demonstrativa* (1697; "Demonstrative Logic") of the Ital-

ian Jesuit Gerolamo Saccheri. Saccheri is better known for his suggestion of the possibility of a non-Euclidean geometry in *Euclides ab Omni Naevo Vindicatus* (1733; "Euclid Cleared of Every Flaw"). Another incisive traditional logic was that of the Dutch philosopher Arnold Geulincx, *Logica fundamentis suis restituta* (1662: "Logic Restored to its Fundamentals"). This work attempted to resurrect the rich detail of scholastic logic, including the theory of *suppositio* and issues of existential import.

LEIBNIZ

With the logical work of the German mathematician, philosopher, and diplomat Gottfried Wilhelm Leibniz, we encounter one of the great triumphs, and tragedies, in the history of logic. He created in the 1680s a symbolic logic that is remarkably similar to George Boole's system of 1847—and Boole is widely regarded as the initiator of mathematical or symbolic logic. But nothing other than vague generalities about Leibniz' goals for logic was published until 1903—well after symbolic logic was in full blossom. Thus one could say that, great though Leibniz' discoveries were, they were virtually without influence in the history of logic. (There remains some slight possibility that Lambert or Boole may have been directly or indirectly influenced by Leibniz' logical system.)

Leibniz' logical research was not entirely symbolic, however, nor was he without influence in the history of (nonsymbolic) logic. Early in his life, Leibniz was strongly interested in the program of Lull, and he wrote the *De arte combinatoria* (1666); this work followed the general Lullian goal of discovering truths by combining concepts into judgments in exhaustive ways and then methodically assessing their truth. Leibniz later developed a goal of devising what he called a "universally characteristic language" (*lingua characteristica universalis*) that would, first, notationally represent concepts by displaying the more basic concepts of which they were composed, and second, naturally represent (in the manner of graphs or pictures, "iconically") the concept in a way that could be easily grasped by readers, no matter what their native tongue. Leibniz studied and was impressed by the method of the Egyptians and Chinese in using picturelike expressions for concepts. The goal of a universal language had already been suggested by Descartes for mathematics as a "universal mathematics"; it had also been discussed extensively by the English philologist George Dalgarno (c. 1626–87) and, for mathematical language and communication, by the French algebraist François Viète (1540–1603). The search for a universal language to replace Latin was seriously taken up again in the late 19th century, first by Giuseppe Peano—whose work on Interlingua, an uninflected form of Latin, was directly inspired by Leibniz' conception—and then with Esperanto. The goal of a logical language also inspired Gottlob Frege, and in the 20th century it prompted the development of the logical language LOGLAN and the computer language PROLOG.

Another and distinct goal Leibniz proposed for logic was a "calculus of reason" (*calculus ratiocinator*). This would naturally first require a symbolism but would then involve explicit manipulations of the symbols according to established rules by which either new truths could be discovered or proposed conclusions could be checked to see if they could indeed be derived from the premises. Reasoning could then take place in the way large sums are done—that is, mechanically or algorithmically—and thus not be subject to individual mistakes and failures of ingenuity. Such derivations could be checked by others or performed by machines, a possibility that Leibniz seriously contemplated. Leibniz' suggestion that machines could be constructed to draw valid inferences or to check the deductions of others was followed up by Charles Babbage, William Stanley Jevons, and Charles Sanders Peirce and his student Allan Marquand in the 19th century, and with wide success on modern computers after World War II.

The symbolic calculus that Leibniz devised seems to have been more of a calculus of reason than a "characteristic" language. It was motivated by his view that most concepts were "composite": they were collections or conjunctions of other more basic concepts. Symbols (letters, lines, or

circles) were then used to stand for concepts and their relationships. This resulted in what is called an "intentional" rather than an "extensional" logic—one whose terms stand for properties or concepts rather than for the things having these properties. Leibniz' basic notion of the truth of a judgment was that the concepts making up the predicate were "included in" the concept of the subject. What Leibniz symbolized as " $A \in B$," or what we might write as " $A = B$ " was that all the concepts making up concept A also are contained in concept B, and vice versa.

Leibniz used two further notions to expand the basic logical calculus. In his notation, " $A \oplus B \in C$ " indicates that the concepts in A and those in B wholly constitute those in C. We might write this as " $A + B = C$ " or " $A \cup B = C$ "—if we keep in mind that A, B, and C stand for concepts or properties, not for individual things. Leibniz also used the juxtaposition of terms in the following way: " $AB \in C$," which we might write as " $A \times B = C$ " or " $A \cap B = C$," signifies in his system that all the concepts in both A and B wholly constitute the concept C.

A universal affirmative judgment, such as "All A's are B's," becomes in Leibniz' notation " $A \in AB$." This equation states that the concepts included in the concepts of both A and B are the same as those in A. A syllogism, "All A's are B's; all B's are C's; therefore all A's are C's," becomes the sequence of equations " $A = AB$; $B = BC$; therefore $A = AC$." This conclusion can be derived from the premises by two simple algebraic substitutions and the associativity of logical multiplication. Leibniz' interpretation of particular and negative statements was more problematic. Although he later seemed to prefer an algebraic, equational symbolic logic, he experimented with many alternative techniques, including graphs.

As with many early symbolic logics, including many developed in the 19th century, Leibniz' system had difficulties with particular and negative statements, and it included little discussion of propositional logic and no formal treatment of quantified relational statements. (Leibniz later became keenly aware of the importance of relations and relational inferences.) Although Leibniz might seem to deserve to be credited with great originality in his symbolic logic—especially in his equational, algebraic logic—it turns out that such insights were relatively common to mathematicians of the 17th and 18th centuries who had a knowledge of traditional syllogistic logic. In 1685 Jakob Bernoulli published a pamphlet on the parallels of logic and algebra and gave some algebraic renderings of categorical statements. Later the symbolic work of Lambert, Ploucquet, Euler, and even Boole—all apparently uninfluenced by Leibniz' or even Bernoulli's work—seems to show the extent to which these ideas were apparent to the best mathematical minds of the day.

THE 18TH AND 19TH CENTURIES

In the 18th century there were three major contributors to the development of formal logic: Ploucquet, Lambert, and Euler, although none went far beyond Leibniz and none influenced subsequent developments in the way that Boole and Frege later did. Leibniz' major goals for logic, such as the development of a "characteristic" language; the parallels among arithmetic, algebra, and syllogistic; and his notion of the truth of a judgment as the concept of the predicate being "included in" the concept of the subject, were carried forward by Christian Wolff but without any significant development of a logic, symbolic or otherwise. The prolific Wolff publicized Leibniz' general views widely and spawned two minor symbolic formulations of logic; that of J.A. Segner in 1740 and that of Joachim Georg Darjes (1714–91) in 1747. Segner used the notation " $B < A$ " to signify, intentionally in the manner of Leibniz, that the concept of B is included in the concept of A (*i.e.*, "All A's are B's").

The work of Gottfried Ploucquet (1716–90) was based on the ideas of Leibniz, although the symbolic calculus Ploucquet developed does not resemble that of Leibniz. The basis of Ploucquet's symbolic logic was the sign ">," which he unfortunately used to indicate that two concepts are disjoint—*i.e.*, having no basic concepts in common; in its propositional interpretation, it is equivalent to what

Intensional
logic

The goal of
a universal
language

Gottfried
Ploucquet

became known in the 20th century as the “Sheffer stroke” function (also known to Peirce) meaning “neither . . . nor.” The universal negative proposition, “No A’s are B’s,” would become “ $A > B$ ” (or, convertibly, “ $B > A$ ”). The equality sign was used to denote conceptual identity, as in Leibniz. Capital letters were used for distributed terms, lowercase ones for undistributed terms. The intersection of concepts was represented by “+”; the multiplication sign (or juxtaposition) stood for the inclusive union of concepts; and a bar over a letter stood for complementation (in the manner of Leibniz). Thus “ \bar{A} ” represented all non-A’s, while “ $\bar{\bar{a}}$ ” meant the same as “some non-A.” Rules of inference were the standard algebraic substitution of identicals along with more complicated implicit rules for manipulating the nonidentities using “ $>$.” Ploucquet was interested in graphic representations of logical relations—using lines, for example. He was also one of the first symbolic logicians to have worried extensively about representing quantification—although his own contrast of distributed and undistributed terms is a clumsy and limited device. Not a mathematician, Ploucquet did not pursue the logical interpretation of inverse operations (e.g., division, square root, and so on) and of binomial expansions; the interpretation of these operations was to plague some algebras of logic and sidetrack substantive development—first in the work of Leibniz and the Bernoullis, then in that of Lambert, Boole, and Schröder. Ploucquet published and promoted his views widely (his publications included an essay on Leibniz’ logic); he influenced his contemporary Lambert and had a still greater influence upon Georg Jonathan von Holland and Christian August Semler.

The greatest 18th-century logician was undoubtedly Johann Heinrich Lambert. Lambert was the first to demonstrate the irrationality of π , and, when asked by Frederick the Great in what field he was most capable, is said to have curtly answered “All.” His own highly articulated philosophy was a more thorough and creative reworking of rationalist ideas from Leibniz and Wolff. His symbolic and formal logic, developed especially in his *Sechs Versuche einer Zeichenkunst in der Vernunftlehre* (1777; “Six Attempts at a Symbolic Method in the Theory of Reason”), was an elegant and notationally efficient calculus, extensively duplicating, apparently unwittingly, sections of Leibniz’ calculus of a century earlier. Like the systems of Leibniz, Ploucquet, and most Germans, it was intensional, using terms to stand for concepts, not individual things. It used an identity sign and the plus sign in the natural algebraic way that one sees in Leibniz and Boole. Five features distinguish it from other systems. First, Lambert was concerned to separate the simpler concepts constituting a more complex concept into the genus and differentia—the broader and narrowing concepts—typical of standard definitions: the symbols for the genus and differentia of a concept were operations on terms, extracting the genus or differentia of a concept. Second, Lambert carefully differentiated among letters for known, undetermined, and genuinely unknown concepts, using different letters from the Latin alphabet; the lack of such distinctions in algebra instruction has probably caused extensive confusion. Third, his disjunction or union operation, “+,” was taken in the exclusive sense—excluding the overlap of two concepts, in distinction to Ploucquet’s inclusive operation, for example. Fourth, Lambert accomplished the expression of quantification such as that in “Every A is B” by writing “ $a = mb$ ”—that is, the known concept a is identical to the concepts in both the known concept b and an indeterminate concept m ; this device is similar enough to Boole’s later use of the letter “ y ” to suggest some possible influence. Finally, Lambert considered briefly the symbolic theorems that would not hold if the concepts were relations, such as “is the father of.” He also introduced a notation for expressing relational notions in terms of single-placed functions: in his system, “ $i = \alpha : : c$ ” indicates that the individual (concept) i is the result of applying a function α to the individual concept c . Although it is not known whether Frege had read Lambert, it is possible that Lambert’s analysis influenced Frege’s analysis of quantified relations, which depends on the notion of a function. Lambert also developed a method of pictorially display-

ing the overlap of the content of concepts with overlapping line segments. Leibniz had experimented with similar techniques. Two-dimensional techniques were popularized by the Swiss mathematician Leonhard Euler in his *Lettres à une princesse d’Allemagne* (1768–74; “Letters to a German Princess”). These techniques and the related Venn diagrams have been especially popular in logic education. In Euler’s method the interior areas of circles represented (intensionally) the more basic concepts making up a concept or property. To display “All A’s are B’s,” Euler drew a circle labeled “A” that was entirely contained within another circle, “B.” Such circles could be manipulated to discover the validity of syllogisms. Euler did not develop this method very far, and it did not constitute a significant logical advance. Leibniz himself had occasionally drawn such illustrations, and they apparently first entered the literature in the *Universalia Euclidea* (1661) of Johann C. Sturm and were more frequently used by Johann C. Lange in 1712. (Vives had employed triangles for similar purposes in 1555.) Euler’s methods were systematically developed by the French mathematician Joseph-Diez Gergonne in 1816–17, although Gergonne retreated from two-dimensional graphs to linear formulas that could be more easily printed and manipulated. For complicated reasons, almost all German formal logic came from the Protestant areas of the German-speaking world.

The German philosophers Immanuel Kant and Georg Wilhelm Friedrich Hegel made enormous contributions to philosophy, but their contributions to formal logic can only be described as minimal or even harmful. Kant refers to logic as a virtually completed artifice in his important *Critique of Pure Reason* (1781). He showed no interest in Leibniz’ goal of a natural, universal, and efficient logical language and no appreciation of symbolic or mathematical formulations. His own lectures on logic, published in 1800 as *Immanuel Kants Logik: ein Handbuch zu Vorlesungen*, and his earlier *The Mistaken Subtlety of the Four Syllogistic Figures* (1762) were minor contributions to the history of logic. Hegel refers early in his massive *Science of Logic* (1812–16) to the centuries of work in logic since Aristotle as a mere preoccupation with “technical manipulations.” He took issue with the claim that one could separate the “logical form” of a judgment from its substance—and thus with the very possibility of logic based on a theory of logical form. When the study of logic blossomed again on German-speaking soil, contributors came from mathematics and the natural sciences.

In the English-speaking world, logic had always been more easily and continuously tolerated, even if it did not so early reach the heights of mathematical sophistication that it had in the German- and French-speaking worlds. Logic textbooks in English appeared in considerable numbers in the 17th and 18th centuries: some were translations, while others were handy, simplified handbooks with some interesting and developed positions, such as John Wallis’ *Institutio Logicae* (1687) and works by Henry Aldrich, Isaac Watts, and the founder of Methodism, John Wesley. Out of this tradition arose Richard Whately’s *Elements of Logic* (1826) and, in the same tradition, John Stuart Mill’s enormously popular *A System of Logic* (1843). Although now largely relegated to a footnote, Whately’s nonsymbolic textbook reformulated many concepts in such a thoughtful and clear way that it is generally (and first by De Morgan) credited with single-handedly bringing about the “rebirth” of English-language logic.

The two most important contributors to British logic in the first half of the 19th century were undoubtedly George Boole and Augustus De Morgan. Their work took place against a more general background of logical work in English by figures such as Whately, George Bentham, Sir William Hamilton, and others. Although Boole cannot be credited with the very first symbolic logic, he was the first major formulator of a symbolic extensional logic that is familiar today as a logic or algebra of classes. (A correspondent of Lambert, Georg von Holland, had experimented with an extensional theory, and in 1839 the English writer Thomas Solly presented an extensional logic in *A Syllabus of Logic*, though not an algebraic one.) Boole published two major works, *The Mathematical*

Leonhard Euler

Johann Heinrich Lambert

Boole and De Morgan

Analysis of Logic in 1847 and *An Investigation of the Laws of Thought* in 1854. It was the first of these two works that had the deeper impact on his contemporaries and on the history of logic. *The Mathematical Analysis of Logic* arose as the result of two broad streams of influence. The first was the English logic-textbook tradition. The second was the rapid growth in the early 19th century of sophisticated discussions of algebra and anticipations of nonstandard algebras. The British mathematicians D.F. Gregory and George Peacock were major figures in this theoretical appreciation of algebra. Such conceptions gradually evolved into “nonstandard” abstract algebras such as quaternions, vectors, linear algebra, and Boolean algebra itself.

Boole used capital letters to stand for the extensions of terms; they are referred to (in 1854) as classes of “things” but should not be understood as modern sets. The universal class or term—which he called simply “the Universe”—was represented by the numeral “1,” and the null class by “0.” The juxtaposition of terms (for example, “AB”) created a term referring to the intersection of two classes or terms. The addition sign signified the non-overlapping union; that is, “ $A + B$ ” referred to the entities in A or in B ; in cases where the extensions of terms A and B overlapped, the expression was held to be “undefined.” For designating a proper subclass of a class, Boole used the notation “ v ,” writing for example “ vA ” to indicate some of the A ’s. Finally, he used subtraction to indicate the removing of terms from classes. For example, “ $1 - x$ ” would indicate what one would obtain by removing the elements of x from the universal class—that is, obtaining the complement of x (relative to the universe, 1).

Basic equations included: $1A = A$, $0A = 0$, $A + 0 = 0$, $A + 1 = 1$ (but only where $A = 0$), $A + B = B + A$, $AB = BA$, $AA = A$ (but not $A + A = A$), $(AB)C = A(BC)$, and the distribution laws, $A(B + C) = AB + AC$ and $A + (BC) = (A + B)(A + C)$. Boole offered a relatively systematic, but not rigorously axiomatic, presentation. For a universal affirmative statement such as “All A ’s are B ’s,” Boole used three alternative notations: $AB = B$ (somewhat in the manner of Leibniz), $A(1 - B) = 0$, or $A = vB$ (the class of A ’s is equal to some proper subclass of the B ’s). The first and second interpretations allowed one to derive syllogisms by algebraic substitution: the latter required manipulation of subclass (“ v ”) symbols.

In contrast to earlier symbolisms, Boole’s was extensively developed, with a thorough exploration of a large number of equations (including binomial-like expansions) and techniques. The formal logic was separately applied to the interpretation of propositional logic, which became an interpretation of the class or term logic—with terms standing for occasions or times rather than for concrete individual things. Following the English textbook tradition, deductive logic is but one half of the subject matter of the book, with inductive logic and probability theory constituting the other half of both his 1847 and 1854 works.

Seen in historical perspective, Boole’s logic was a remarkably smooth bend of the new “algebraic” perspective and the English-logic textbook tradition. His 1847 work begins with a slogan that could have served as the motto of abstract algebra: “. . . the validity of the processes of analysis does not depend upon the interpretation of the symbols which are employed, but solely upon the laws of combination.”

Modifications to Boole’s system were swift in coming: in the 1860s Peirce and Jevons both proposed replacing Boole’s “+” with a simple inclusive union or summation: the expression “ $A + B$ ” was to be interpreted as designating the class of things in A , in B , or in both. This results in accepting the equation “ $1 + 1 = 1$,” which is certainly not true of the ordinary numerical algebra and at which Boole apparently balked.

Interestingly, one defect in Boole’s theory, its failure to detail relational inferences, was dealt with almost simultaneously with the publication of his first major work. In 1847 Augustus De Morgan published his *Formal Logic; or, the Calculus of Inference, Necessary and Probable*. Unlike Boole and most other logicians in the United Kingdom, De Morgan knew the medieval theory of logic and semantics and also knew the Continental, Leibnizian

symbolic tradition of Lambert, Ploucquet, and Gergonne. The symbolic system that De Morgan introduced in his work and used in subsequent publications is, however, clumsy and does not show the appreciation of abstract algebras that Boole’s did. De Morgan did introduce the enormously influential notion of a possibly arbitrary and stipulated “universe of discourse” that was used by later Booleans. (Boole’s original universe referred simply to “all things.”) This view influenced 20th-century logical semantics. De Morgan contrasted uppercase and lowercase letters: a capital letter represented a class of individuals, while a lowercase letter represented its complement relative to the universe of discourse, a convention Boole might have expressed by writing “ $x = (1 - X)$ ”; this stipulation results in the general principle: $xX = 0$. A period indicated a (propositional) negation, and the parentheses (“and”) indicated, respectively, distributed (if the parentheses faces toward the nearby term) and undistributed terms. Thus De Morgan would write “All A ’s are B ’s” as “ $A)B$ ” and “Some A ’s are B ’s” as “ $A ()B$.” These distinctions parallel Boole’s account of distribution (quantification) in “ $A = vB$ ” (where A is distributed but B is not) and “ $vA = B$ ” (where both terms are distributed). Although his entire system was developed with wit, consistency, and brilliance, it is remarkable that De Morgan never saw the inferiority of his notation to almost all available symbolisms.

De Morgan’s other essays on logic were published in a series of papers from 1846 to 1862 (and an unpublished essay of 1868) entitled simply “On the Syllogism.” The first series of four papers found its way into the middle of the *Formal Logic* of 1847. The second series, published in 1850, is of considerable significance in the history of logic, for it marks the first extensive discussion of quantified relations since late medieval logic and Jung’s massive *Logica hamburgensis* of 1638. In fact, De Morgan made the point, later to be exhaustively repeated by Peirce and implicitly endorsed by Frege, that relational inferences are the core of mathematical inference and scientific reasoning of all sorts; relational inferences are thus not just one type of reasoning but rather are the most important type of deductive reasoning. Often attributed to De Morgan—not precisely correctly but in the right spirit—was the observation that all of Aristotelian logic was helpless to show the validity of the inference, “All horses are animals; therefore, every head of a horse is the head of an animal.” The title of this series of papers, De Morgan’s devotion to the history of logic, his reluctance to mathematize logic in any serious way, and even his clumsy notation—apparently designed to represent as well as possible the traditional theory of the syllogism—show De Morgan to be a deeply traditional logician.

Charles Sanders Peirce, the son of the Harvard mathematics professor and discoverer of linear algebra Benjamin Peirce, was the first significant American figure in logic. Peirce had read the work of Aristotle, Whately, Kant, and Boole as well as medieval works and was influenced by his father’s sophisticated conceptions of algebra and mathematics. Peirce’s first published contribution to logic was his improvement in 1867 of Boole’s system. Although Peirce never published a book on logic (he did edit a collection of papers by himself and his students, the *Studies in Logic* of 1883), he was the author of an important article in 1870, whose abbreviated title was “On the Notation of Relatives,” and of a series of articles in the 1880s on logic and mathematics; these were all published in American mathematics journals.

It is relatively easy to describe Peirce’s main approach to logic, at least in his earlier work: it was a refinement of Boole’s algebra of logic and, especially, the development of techniques for handling relations within that algebra. In a phrase, Peirce sought a blend of Boole (on the algebra of logic) and De Morgan (on quantified relational inferences). Described in this way, however, it is easy to underestimate the originality and creativity (even idiosyncrasy) of Peirce. Although committed to the broadly “algebraic” tradition of Boole and his father, Peirce quickly moved away from the equational style of Boole and from efforts to mimic numerical algebra. In particular, he argued that

Charles
Sanders
Peirce

a transitive and asymmetric logical relation of inclusion, for which he used the symbol " \prec ," was more useful than equations; the importance of such a basic, transitive relation was first stressed by De Morgan, and much of Peirce's work can be seen as an exploration of the formal, abstract properties of this distinctively logical relation. He used it to express class inclusion, the "if . . . then" connective of propositional logic, and even the relation between the premises and conclusion of an argument ("illation"). Furthermore, Peirce slowly abandoned the strictly substitutional character of algebraic terms and increasingly used notation that resembled modern quantifiers (see above *Formal logic: The predicate calculus*). Quantifiers were briefly introduced in 1870 and were used extensively in the papers of the 1880s. They were borrowed by Schröder for his extremely influential treatise on the algebra of logic and were later adopted by Peano from Schröder; thus in all probability they are the source of the notation for quantifiers now widely used. In his earlier works, Peirce might have written " $A \prec B$ " to express the universal statement "All A 's are B 's"; however, he often wrote this as " $\Pi_i A_i \prec \Pi_i B_i$ " (the class of all the i 's that are A is included in the class of all the i 's that are B) or, still later and interpreted in the modern way, as "For all i 's, if i is A , then i is B ." Peirce and Schröder were never clear about whether they thought these quantifiers and variables were necessary for the expression of certain statements (as opposed to using strictly algebraic formulas), and Frege did not address this vital issue either; the Boolean algebra without quantifiers, even with extensions for relations that Peirce introduced, was demonstrated to be inadequate only in the mid-20th century by Alfred Tarski and others.

Algebra
of relations
and
existential
graphs

Peirce developed this symbolism extensively for relations. His earlier work was based on versions of multiplication and addition for relations—called relative multiplication and addition—so that Boolean laws still held. Both Peirce's conception of the purposes of logic and the details of his symbolism and logical rules were enormously complicated by highly developed and unusual philosophical views, by elaborate theories of mind and thought, and by his theory of mental and visual signs (semiotics). He argued that all reasoning was "diagrammatic" but that some diagrams were better (more iconic) than others if they more accurately represented the structure of our thoughts. His earlier works seem to be more in the tradition of developing a calculus of reason that would make reasoning quicker and better and permit one to validate others' reasoning more accurately and efficiently. His later views, however, seem to be more in the direction of developing a "characteristic" language. In the late 1880s and 1890s Peirce developed a far more extensively iconic system of logical representation, his existential graphs. This work was, however, not published in his lifetime and was little recognized until the 1960s.

Peirce did not play a major role in the important debates at the end of the 19th century on the relationship of logic and mathematics and on set theory. In fact, in responding to an obviously quick reading of Russell's restatements of Frege's position that mathematics could be derived from logic, Peirce countered that logic was properly seen as a branch of mathematics, not vice versa. He had no influential students: the brilliant O.H. Mitchell died at an early age, and Christine Ladd Franklin never adapted to the newer symbolic tradition of Peano, Frege, and Russell. On the other hand, Peano and especially Schröder had read Peirce's work carefully and adopted much of his notation and his doctrine of the importance of relations (although they were less fervent than De Morgan and Peirce). Peano and Schröder, using much of Peirce's notation, had an enormous influence into the 20th century.

In Germany, the older formal and symbolic logical tradition was barely kept alive by figures such as Salomon Maimon, Semler, August Detlev Twisten, and Moritz Wilhelm Drobisch. The German mathematician and philologist Hermann Günther Grassmann published in 1844 his *Ausdehnungslehre* ("The Theory of Extension"), in which he used a novel and difficult notation to explore quantities ("extensions") of all sorts—logical extension and intension, numerical, spatial, temporal, and so on. Grassmann's

notion of extension is very similar to the use of the broad term "quantity" (and the phrase "logic of quantity") that is seen in the works of George Bentham and Sir William Hamilton from the same period in the United Kingdom; it is from this English-language tradition that the terms, still in use, of logical "quantification" and "quantifiers" derive. Grassmann's work influenced Robert Grassmann's *Die Begriffslehre oder Logik* (1872; "The Theory of Concepts or Logic"), Schröder, and Peano. The stage for a rebirth of German formal logic was further set by Friedrich Adolf Trendelenburg's works, published in the 1860s and '70s, on Aristotle's and Leibniz' logic and on the relationship of mathematics and philosophy. Alois Riehl's much-read article "Die englische Logik der Gegenwart (1876; "Contemporary English Logic") introduced German speakers to the works of Boole, De Morgan, and Jevons.

In 1879 the young German mathematician Gottlob Frege—whose mathematical specialty, like Boole's, had actually been calculus—published perhaps the finest single book on symbolic logic in the 19th century, *Begriffsschrift* ("Conceptual Notation"). The title was taken from Trendelenburg's translation of Leibniz' notion of a characteristic language. Frege's small volume is a rigorous presentation of what would now be called the first-order predicate logic. It contains a careful use of quantifiers and predicates (although predicates are described as functions, suggestive of the technique of Lambert). It shows no trace of the influence of Boole and little trace of the older German tradition of symbolic logic. One might surmise that Frege was familiar with Trendelenburg's discussion of Leibniz, had probably encountered works by Drobisch and Hermann Grassmann, and possibly had a passing familiarity with the works of Boole and Lambert, but was otherwise ignorant of the history of logic. He later characterized his system as inspired by Leibniz' goal of a characteristic language but not of a calculus of reason. Frege's notation was unique and problematically two-dimensional; this alone caused it to be little read.

Gottlob
Frege

Frege was well aware of the importance of functions in mathematics, and these form the basis of his notation for predicates; he never showed an awareness of the work of De Morgan and Peirce on relations or of older medieval treatments. The work was reviewed (by Schröder, among others), but never very positively, and the reviews always chided him for his failure to acknowledge the Boolean and older German symbolic tradition; reviews written by philosophers chided him for various sins against reigning idealist dogmas. Frege stubbornly ignored the critiques of his notation and persisted in publishing all his later works using it, including his little-read magnum opus, *Grundgesetze der Arithmetik* (1893–1903; *The Basic Laws of Arithmetic*).

His first writings after the *Begriffsschrift* were bitter attacks on Boolean methods (showing no awareness of the improvements by Peirce, Jevons, Schröder, and others) and a defense of his own system. His main complaint against Boole was the artificiality of mimicking notation better suited for numerical analysis rather than developing a notation for logical analysis alone. This work was followed by the *Die Grundlagen der Arithmetik* (1884; *The Foundations of Arithmetic*) and then by a series of extremely important papers on precise mathematical and logical topics. After 1879 Frege carefully developed his position that all of mathematics could be derived from, or reduced to, basic "logical" laws—a position later to be known as logicism in the philosophy of mathematics. His view paralleled similar ideas about the reducibility of mathematics to set theory from roughly the same time—although Frege always stressed that his was an intensional logic of concepts, not of extensions and classes. His views are often marked by hostility to British extensional logic and to the general English-speaking tendencies toward nominalism and empiricism that he found in authors such as J.S. Mill. Frege's work was much admired in the period 1900–10 by Bertrand Russell who promoted Frege's logicist research program—first in the *Introduction to Mathematical Logic* (1903), and then with Alfred North Whitehead, in *Principia Mathematica* (1910–13)—but who used a Peirce-Schröder-Peano system of notation rather than

The logicist
position

Frege's; Russell's development of relations and functions was very similar to Schröder's and Peirce's. Nevertheless, Russell's formulation of what is now called the "set-theoretic" paradoxes was taken by Frege himself, perhaps too readily, as a shattering blow to his goal of founding mathematics and science in an intensional, "conceptual" logic. Almost all progress in symbolic logic in the first half of the 20th century was accomplished using set theories and extensional logics and thus mainly relied upon work by Peirce, Schröder, Peano, and Georg Cantor. Frege's care and rigour were, however, admired by many German logicians and mathematicians, including David Hilbert and Ludwig Wittgenstein. Although he did not formulate his theories in an axiomatic form, Frege's derivations were so careful and painstaking that he is sometimes regarded as a founder of this axiomatic tradition in logic. Since the 1960s Frege's works have been translated extensively into English and reprinted in German, and they have had an enormous impact on a new generation of mathematical and philosophical logicians.

German symbolic logic (in a broad sense) was cultivated by two other major figures in the 19th century. The tradition of Hermann Grassmann was continued by the German mathematician and algebraist Ernst Schröder. His first work, *Der Operations-kreis des Logikkalküls* (1877; "The Circle of Operations of the Logical Calculus"), was an equational algebraic logic influenced by Boole and Grassmann but presented in an especially clear, concise, and careful manner; it was, however, intensional in that letters stand for concepts, not classes or things. Although Jevons and Frege complained of what they saw as the "mysterious" relationship between numerical algebra and logic in Boole, Schröder announced with great clarity: "There is certainly a contrast of the objects of the two operations. They are totally different. In arithmetic, letters are numbers, but here, they are arbitrary concepts." He also used the phrase "mathematical logic." Schröder's main work was his three-volume *Vorlesungen über die Algebra der Logik* (1890–1905; "Lectures on the Algebra of Logic"). This is an extensive and sometimes original presentation of all that was known about the algebra of logic circa 1890, together with derivations of thousands of theorems and an extensive bibliography of the history of logic. It is an extensional logic with a special sign for inclusion " \subset " (paralleling Peirce's " \prec "), an inclusive notion of class union, and the usual Boolean operations and rules.

The first volume is devoted to the basic theory of an extensional theory of classes (which Schröder called *Gebiete*, logical "domains," a term that is somewhat suggestive of Grassmann's "extensions"). Schröder was especially interested in formal features of the resulting calculus, such as the property he called "dualism" (carried over from his 1877 work): any theorem remains valid if the addition and multiplication, as well as 0 and 1, are switched—for example, $A \bar{A} = 0$, $A + \bar{A} = 1$, and the pair of De Morgan laws. The second volume is a discussion of propositional logic, with propositions taken to refer to domains of times in the manner of Boole's *Laws of Thought* but using the same calculus. Schröder, unlike Boole and Peirce, distinguished between the universes for the separate cases of the class and propositional logics, using respectively 1 and $\bar{1}$. The third volume contains Schröder's masterful but leisurely development of the logic of relations, borrowing heavily from Peirce's work. In the first decades of the 20th century, Schröder's volumes were the only major works in German on symbolic logic other than Frege's, and they had an enormous influence on important figures writing in German, such as Thoralf Albert Skolem, Leopold Löwenheim, Julius König, Hilbert, and Tarski. (Frege's influence was felt mainly through Russell and Whitehead's *Principia Mathematica*, but this tradition had a rather minor impact on 20th-century German logic.) Although it was an extensional logic more in the English tradition, Schröder's logic exhibited the German tendency of focusing exclusively upon deductive logic; it was a legacy of the English textbook tradition always to cover inductive logic in addition, and this trait survived in (and often cluttered) the works of Boole, De Morgan, Venn, and Peirce.

A development in Germany originally completely dis-

tingent from logic but later to merge with it was Georg Cantor's development of set theory. In work originating from discussions on the foundations of the infinitesimal and derivative calculus by Baron Augustin-Louis Cauchy and Karl Weierstrauss, Cantor and Richard Dedekind developed methods of dealing with the large, and in fact infinite, sets of the integers and points on the real number line. Although the Booleans had used the notion of a class, they rarely developed tools for dealing with infinite classes, and no one systematically considered the possibility of classes whose elements were themselves classes, which is a crucial feature of Cantorian set theory. The conception of "real" or "closed" infinities of things, as opposed to infinite possibilities, was a medieval problem that had also troubled 19th-century German mathematicians, especially the great Carl Friedrich Gauss. The Bohemian mathematician and priest Bernhard Bolzano emphasized the difficulties posed by infinities in his *Paradoxien des Unendlichen* (1851; "Paradoxes of the Infinite"); in 1837 he had written an anti-Kantian and pro-Leibnizian nonsymbolic logic that was later widely studied. First Dedekind, then Cantor used Bolzano's tool of measuring sets by one-to-one mappings; using this technique, Dedekind gave in *Was sind und was sollen die Zahlen?* (1888; "What Are and Should Be the Numbers?") a precise definition of an infinite set. A set is infinite if and only if the whole set can be put into one-to-one correspondence with a proper part of the set. (De Morgan and Peirce had earlier given quite different but technically correct characterizations of infinite domains; these were not especially useful in set theory and went unnoticed in the German mathematical world.)

Although Cantor developed the basic outlines of a set theory, especially in his treatment of infinite sets and the real number line, he did not worry about rigorous foundations for such a theory—thus, for example, he did not give axioms of set theory—nor about the precise conditions governing the concept of a set and the formation of sets. Although there are some hints in Cantor's writing of an awareness of problems in this area (such as hints of what later came to be known as the class/set distinction), these difficulties were forcefully posed by the paradoxes of Russell and the Italian mathematician Cesare Burali-Forti and were first overcome in what has come to be known as Zermelo-Fraenkel set theory.

French logic was ably, though not originally, represented in this period by Louis Liard and Louis Couturat. Couturat's *L'Algèbre de la logique* (1905; *The Algebra of Logic*) and *De l'Infini mathématique* (1896; "On Mathematical Infinity") were important summaries of German and English research on symbolic logic, while his book on Leibniz' logic (1901) and an edition of Leibniz' previously unpublished writings on logic (1903) were very important events in the study of the history of logic. In Russia V.V. Bobylin (1886) and Platon Sergeevich Poretsky (1884) initiated a school of algebraic logic. In the United Kingdom a vast amount of work on formal and symbolic logic was published in the best philosophical journals from 1870 until 1910. This includes work by William Stanley Jevons, whose intensional logic is unusual in the English-language tradition; John Venn, who was notable for his (extensional) diagrams of class relationships but who retained Boole's noninclusive class union operator; Hugh MacColl; Alexander Bain; Sophie Bryant; Emily Elizabeth Constance Jones; Arthur Thomas Shearman; Lewis Carroll (Charles Lutwidge Dodgson); and Whitchhead, whose *A Treatise on Universal Algebra* (1898) was the last major English logical work in the algebraic tradition. Little of this work influenced Russell's conception, which was soon to sweep through English-language logic; Russell was more influenced by Frege, Peano, and Schröder. The older nonsymbolic syllogistic tradition was represented in major English universities well into the 20th century by John Cook Wilson, William Ernest Johnson, Lizzie Susan Stebbing, and Horace William Brindley Joseph and in the United States by Ralph Eaton, James Edwin Creighton, Charles West Churchman, and Daniel Sommer Robinson.

The Italian mathematician Giuseppe Peano's contributions represent a more extensive impetus to the new, nonalgebraic logic. He had a direct influence on the no-

Ernst
Schröder

Cantorian
set theory

British
logicians

tation of later symbolic logic that exceeded that of Frege and Peirce. His early works (such as the logical section of the *Calcolo geometrico secondo l'Ausdehnungslehre di H. Grassman* [1888; "Calculus of Geometry According to the Theory of Extension of H. Grassmann"]) were squarely in the algebraic tradition of Boole, Grassmann, Peirce, and Schröder. Writing in the 1890s in his own journal, *Revista di matematica*, with a growing appreciation of the use of quantifiers in the first and third volumes of Schröder's *Vorlesungen*, Peano evolved a notation for quantifiers. This notation, along with Peano's use of the Greek letter epsilon, ϵ , to denote membership in a set, was adopted by Russell and Whitehead and used in later logic and set theory. Although Peano himself was not interested in the logicist program of Frege and Russell, his five postulates governing the structure of the natural numbers (now known as the Peano Postulates), with similar ideas in the work of Peirce and Dedekind, came to be regarded as the crucial link between logic and mathematics. It was widely thought that all mathematics could be derived from the theory for the natural numbers; if the Peano postulates could be derived from logic, or from logic including set theory, its feasibility would have been demonstrated. Simultaneously with his work in logic, Peano wrote many articles on universal languages and on the features of an ideal notation in mathematics and logic—all explicitly inspired by Leibniz.

The Paris international congresses

Logic in the 19th century culminated grandly with the First International Congress of Philosophy and the Second International Congress of Mathematics held consecutively in Paris in August 1900. The overlap between the two congresses was extensive and fortunate for the future of logic and philosophy. Peano, Alessandro Padoa, Burali-Forti, Schröder, Cantor, Dedekind, Frege, Felix Klein, Ladd Franklin (Peirce's student), Coutourat, and Henri Poincaré were on the organizing committee of the Philosophical Congress; for the subsequent development of logic, Bertrand Russell was perhaps its most important attendee. The influence of algebraic logic was already ebbing, and the importance of nonalgebraic symbolic logics, of axiomatizations, and of logic (and set theory) as a foundation for mathematics were ascendant. Until the congresses of 1900 and the work of Russell and Hilbert, mathematical logic lacked full academic legitimacy. None of the 19th-century logicians had achieved major positions at first-rank universities: Peirce never obtained a permanent university position, Dedekind was a high-school teacher, and Frege and Cantor remained at provincial universities. The mathematicians stayed for the Congress of Mathematics, and it was here that David Hilbert gave his presentation of the 23 most significant unsolved problems of mathematics—several of which were foundational issues in mathematics and logic that were to dominate logical research during the first half of the 20th century.

20th-century logic

In 1900 logic was poised on the brink of the most active period in its history. The late 19th-century work of Frege, Peano, and Cantor, as well as Peirce's and Schröder's extensions of Boole's insights, had broken new ground, raised considerable interest, established international lines of communication, and formed a new alliance between logic and mathematics. Five projects internal to late 19th-century logic coalesced in the early 20th century; especially in works such as Russell and Whitehead's *Principia Mathematica*. These were the development of a consistent set or property theory (originating in the work of Cantor and Frege), the application of the axiomatic method (including non-symbolically), the development of quantificational logic, and the use of logic to understand mathematical objects and the nature of mathematical proof. The five projects were unified by a general effort to use symbolic techniques, sometimes called mathematical, or formal, techniques. Logic became increasingly "mathematical," then, in two senses. First, it attempted to use symbolic methods like those that had come to dominate mathematics. Second, an often dominant purpose of logic came to be its use as a tool for understanding the na-

ture of mathematics—such as in defining mathematical concepts, precisely characterizing mathematical systems, or describing the nature of ideal mathematical proof. (See MATHEMATICS, THE HISTORY OF: *Mathematics in the 19th and 20th centuries*, and MATHEMATICS, THE FOUNDATIONS OF.)

RUSSELL AND WHITEHEAD'S PRINCIPIA MATHEMATICA

The three-volume *Principia Mathematica* (1910–13) was optimistically named after the *Philosophiae naturalis principia mathematica* of another hugely important Cambridge thinker, Isaac Newton. Like Newton's *Principia*, it was imbued with an optimism about the application of mathematical techniques, this time not to physics but to logic and to mathematics itself—what the first sentence of their preface calls "the mathematical treatment of the principles of mathematics." It was intended by Russell and Whitehead both as a summary of then-recent work in logic (especially by Frege, Cantor and Peano) and as a ground-breaking, large-scale treatise systematically developing mathematical logic and deriving basic mathematical principles from the principles of logic alone.

The *Principia* was the natural outcome of Russell's earlier polemical book, *The Principles of Mathematics* (published in 1903 but largely written in 1900), and his views were later summarized in *Introduction to Mathematical Philosophy* (1919). Whitehead's *A Treatise on Universal Algebra* (1898) was more in the algebraic tradition of Boole, Peirce, and Schröder, but there is a sense in which *Principia Mathematica* became the second volume both of it and of Russell's *Principles*.

The main idea in the *Principia* is the view, taken from Frege, that all of mathematics could be derived from the principles of logic alone. This view later came to be known as logicism and was one of the principal philosophies of mathematics in the early 20th century. Number theory, the core of mathematics, was organized around the Peano postulates, stated in works by Peano of 1889 and 1895 (and anticipated by similar but less influential theories of Peirce and Dedekind). These postulates state and organize the fundamental laws of "natural" (integral, positive) numbers, and thus of all of mathematics:

The Peano postulates

1. 0 is a number.
2. The successor of any number is also a number.
3. No two distinct numbers have the same successor.
4. 0 is not the successor of any number.
5. If any property is possessed by 0 and also by the successor of any number having the property, then all numbers have that property.

If some entities satisfying these conditions could be derived or constructed in logic, it would have been shown that mathematics was (or at least could be) founded in pure logic, requiring no additional assumptions.

Although his language actually used the intensional and second-order language of functions and properties, Frege had claimed to have accomplished precisely this, identifying 0 with the empty set, 1 with the set of all single-membered sets (singletons), 2 with the set of all dual-membered sets (doubletons), and so on. These sets of equinumerous sets were then what numbers really were. Unfortunately, Russell showed through his famous paradox that the theory is inconsistent and, hence, that any statement at all can be derived in Frege's system, not merely desired logical truths, the Peano postulates, and what follows from them. Russell, in a famous letter to Frege, asked him to consider "the set of all those sets not members of themselves." Paradox follows if one assumes such a set is empty, or is not empty. After meditating on this paradox and a great many other paradoxes devised by Burali-Forti, George Godfrey Berry, and others, Russell and Whitehead concluded that the main difficulty lies in allowing the construction of entities that contain a "vicious circle"—*i.e.*, entities that are used in the construction or definition of themselves.

Russell and Whitehead sought to rule out this possibility while at the same time allowing a great many of the operations that Frege had deemed desirable. The result was the theory of types: all sets and other entities have a logical "type," and sets are always constructed from spec-

ifying members with lower types. (F.P. Ramsay offered a criticism that was subsequently accommodated in later editions of *Principia Mathematica*; as modified, the theory came to be known as the “ramified” theory of types.) Consequently, to speak of sets that are, or are not, “members of themselves” is simply to violate this rule governing the specification of sets. There is some evidence that Cantor had been aware of the difficulties created when there is no such restriction (he permitted large collective entities that do not obey the usual rules for sets), and a parallel intuition concerning the pitfalls of certain operations was independently followed by Ernst Zermelo in the development of his set theory.

In addition to its notation (much of it borrowed from Peano), its masterful development of logical systems for propositional and predicate logic, and its overcoming of difficulties that had beset earlier logical theories and logistic conceptions, the *Principia* offered discussions of functions, definite descriptions, truth, and logical laws that had a deep influence on discussions in analytical philosophy and logic throughout the 20th century. What is perhaps missing is any hesitation or perplexity about the limits of logic: whether this logic is, for example, provably consistent, complete, or decidable, or whether there are concepts expressible in natural languages but not in this logical notation. This is somewhat odd, given the well-known list of problems posed by Hilbert in 1900 that came to animate 20th-century logic, especially German logic. The *Principia* is a work of confidence and mastery and not of open problems and possible difficulties and shortcomings; it is a work closer to the naive progressive elements of the *Jahrhundertwende* than to the agonizing *fin de siècle*.

20TH-CENTURY SET THEORY

Independently of Russell and Whitehead’s work, and more narrowly in the German mathematical tradition of Dedekind and Cantor, in 1908 Ernst Zermelo described axioms of set theory that, slightly modified, came to be standard in the 20th century. The type theory of the *Principia Mathematica* has, by contrast, gradually faded in influence. Like that of Russell and Whitehead, Zermelo’s system avoids the paradoxes inherent in Frege’s and Cantor’s systems by imposing certain restrictions on what may be a set.

Zermelo’s axioms are:

1. Axiom of extensionality. If two sets have the same members, then they are identical.
2. Axiom of elementary sets. There exists a set with no members, the null or empty set. For any two members of a set, there exist (singleton) sets containing only those members, as well as a (doubleton) set containing only those members.
3. Axiom of separation. For any well-formed property and any set S , there is a set, S' , containing all and only the members of S having this property. That is, already-existing sets can be partitioned or separated into parts by certain properties.
4. Power set axiom. If S is a set, then there exists a set, S' , which contains all and only the subsets of S .
5. Union axiom. If S is a set, then there is a set containing all and only the members of the sets in S .
6. Axiom of choice. (Discussed below.)
7. Axiom of infinity. There exists at least one set that contains an infinite number of members.

With the exception of axiom 2, all these axioms allow new sets to be constructed from already-constructed sets by carefully constrained operations. This method embodies what has come to be known as the “iterative” conception of a set. This list of axioms was eventually modified by Zermelo, and by Abraham Fraenkel, and the result is widely known as Zermelo-Fraenkel set theory, or ZF for short. (See the article SET THEORY.)

Axiomatized Zermelo-Frankel set theory is almost always what mathematicians and logicians now mean by “set theory.” The system was later modified by John von Neumann and others with the addition of a “foundation axiom” explicitly prohibiting (among others) sets that contain themselves as members. The system was further modified for technical reasons by von Neumann, Paul

Bernays, and Kurt Gödel in the 1920s and '30s, and the result is called von Neumann-Bernays-Gödel set theory, or NBG. A more distinct alternative was proposed by the American logician Willard Van Orman Quine and is called New Foundations (NF; from 1936–37). Quine’s system is not widely used, however, and there have been recurrent suspicions that it is inconsistent. Other set theories have been proposed, but most of them, such as relevant, fuzzy, or multivalued set theories, differ from ZF in having different underlying logics. ZF was soon shown to be capable of deriving the Peano postulates by several alternative methods—for example, by identifying the natural numbers with certain sets, such as 0 with the empty set, Λ , 1 with the singleton empty set $\{\Lambda\}$, and so on. The crucial mathematical notions of relation and function were defined as certain sets of ordered pairs, and ordered pairs were defined strictly within set theory using suggestions first by the American mathematician and cyberneticist Norbert Wiener, and then by the Polish logician Kasimierz Kuratowski and the Norwegian logician Thoralf Skolem. With these proposals, the need for basic notions of function or relation (and generally, of order) that had been proposed by Frege, Peirce, and Schröder, disappeared. Zermelo and other early set theorists (obviously influenced by Hilbert’s list of open problems) were concerned with a number of issues about the properties of the whole system: Was ZF consistent? Was its consistency provable? Were the axioms independent of one another? Were there other desirable axioms that should be added? Particularly problematic was the status of the axiom of choice.

The axiom of choice states (in Zermelo’s first version) that, given any set of disjoint (nonoverlapping) sets, a set can be formed with one and only one element from each of these disjoint sets. The issue is whether elements can be “chosen” or selected from sets; the problem is acute only when infinite sets are permitted or when numerous nonidentical memberless sets (similar to the empty set), which Zermelo called *Urelementen* (literally “primitive” or “original” elements), are permitted. The axiom of choice has a large number of formulations that are logically equivalent to it, some quite surprisingly so: these include a well-ordering axiom (that the elements of any set can be put in a certain order). Early perplexity in set theory centred on whether the axiom of choice is consistent with the other axioms and whether or not it is independent of them. While clearly desirable, the axiom of choice has the nonintuitive character of a postulate, rather than being self-evident. The first question is whether the addition of the axiom of choice to a system of axiomatic set theory makes the resulting system inconsistent if it was not so previously. The second question is whether the axiom of choice can be derived from the other axioms or whether its inclusion really adds anything to the system—*i.e.*, whether every useful implication of it could be derived without it. The consistency of the axiom of choice with the other axioms of set theory (specifically in NBG set theory) was shown by Kurt Gödel in 1940. The independence of the axiom of choice from the other axioms was shown, trivially, for set theories with *Urelementen* very early; the independence of the axiom of choice in NBG or ZF set theories was one of the major outstanding problems in 20th-century mathematical logic until Paul Cohen showed in 1963 that the axiom of choice was indeed independent of the other standard axioms for set theory.

Another early outstanding issue in axiomatic set theory was whether what came to be known as the “continuum hypothesis” was consistent with the other axioms of ZF, and whether it was independent of them. The continuum hypothesis states that between \aleph_1 (aleph-null; the “smallest” infinite cardinality, on the order of the integers) and its power set, \aleph_2 (a cardinality usually identified with the continuum of points on a real number line), as well as between other integral alephs, there is no intermediate cardinality—no $\aleph_{1.5}$, so to speak. This is the first of Hilbert’s 1900 list of 23 open problems in mathematics and its foundations. The second problem is the consistency and independence of the Peano postulates and any alternative general axioms for mathematics. In his work

Zermelo’s
axioms of
set theory

The axiom
of choice

of 1938–40, Kurt Gödel had shown that the continuum hypothesis—that there are no intermediate cardinalities—was consistent with the other axioms of set theory. In 1963, employing techniques similar to those that he used for showing the independence of the axiom of choice, Paul Cohen showed the independence of the continuum hypothesis. Since 1963 a number of alternative and less difficult methods of showing these independence results have emerged.

A third, but less conceptually vital, area of research in set theory has been in the precise form of axioms of infinity. It became evident that there are a variety of “stronger” axioms of infinity that can be added to ZF: these declare the existence of infinite sets with cardinalities beyond all the integral \aleph cardinalities. With the results of Gödel from 1931, which have implications for the completeness and consistency of set theory (and are discussed below), and with the independence results of Cohen from 1963, basic questions concerning standard set theories (ZF and NBG) are considered to have been answered, even if the answers are somewhat unsatisfactory. The questions that have lingered about set theory, now a very well understood formal system, have centred on philosophical issues of whether numbers and mathematical operations are really “just” sets and set-theoretic operations, or whether one can usefully understand mathematics and the world in other than set-theoretic terms.

Alternatives to set theory

Various substantive alternatives to set theory have been proposed. One is the part-whole calculus, or “calculus of individuals,” also called mereology, of Stanisław Leśniewski (1916, 1927–31). This theory rejects the hierarchy of sets, sets of sets, and so on, that emerge in set theory through the member-of relation (as defined by the power set axiom) and instead proposes a part-whole relationship. It obeys rules like those for the subset relationship in set theory. Its inspiration seems to have been the earlier Boolean theory of classes (especially as described by Schröder), as well as the work of the German philosopher Edmund Husserl and his followers on conceptualization in everyday thought (“Phenomenology”) of collections. This work was developed by Henry Leonard and Nelson Goodman in the United States in the mid-20th century. It has continued to attract philosophers of logic and mathematics who are nominalists, who suspect set theory of being inherently Platonistic, or who are otherwise suspicious of the complex entities proposed by, and the complicated assumptions needed for, set theory. Although some interesting proposals have been made, it does not appear that the part-whole calculus is capable of grounding mathematics, or at least of doing so in as straightforward a manner as does ZF. A much different approach to logical foundations for mathematics is to be seen in the category theory of Saunders MacLane and others. The category theory proposes that mathematics is based on highly abstract formal objects: categories (“topoi,” singular: “topos”) that are neither sets nor properties. In set theory there is a distinction between the objects of the theory—sets—and what one does to them: intersecting them, unioning them, and so forth. In category theory, this distinction between objects and operations on them (transformations, or morphisms) disappears. In the latter part of the 20th century, interest has also arisen in the logic of collective entities other than sets, classes, or classes of individuals: this includes theories of heaps and aggregates and theories for mass terms—such as water or butter—that are not conceptualized as formed from distinct individuals. The goal has been to give formal theories for collective or quantitative terms used in natural language.

LOGIC AND PHILOSOPHIES OF MATHEMATICS

Philosophies of mathematics are more extensively discussed in the article MATHEMATICS, THE FOUNDATIONS OF; the major schools are mentioned here briefly. An outgrowth of the theory of Russell and Whitehead, and of most modern set theories, was a better articulation of a philosophy of mathematics known as logicism: that operations and objects spoken about in mathematics are really purely logical constructions. This has focused increased attention on what exactly “pure” logic is and whether, for

example, set theory is really logic in a narrow sense. There seems little doubt that set theory is not “just” logic in the way in which, for example, Frege viewed logic—*i.e.*, as a formal theory of functions and properties. Because set theory engenders a large number of interestingly distinct kinds of nonphysical, nonperceived abstract objects, it has also been regarded by some philosophers and logicians as suspiciously (or endearingly) Platonistic. Others, such as Quine, have “pragmatically” endorsed set theory as a convenient way—perhaps the only such way—of organizing the whole world around us, especially if this world contains the richness of transfinite mathematics.

For most of the first half of the 20th century, new work in logic saw logic’s goal as being primarily to provide a foundation for, or at least to play an organizing role in, mathematics. Even for those researchers who did not endorse the logicist program, logic’s goal was closely allied with techniques and goals in mathematics, such as giving an account of formal systems (formalism) or of the ideal nature of nonempirical proof and demonstration. (Interest in the logicist and formalist program waned after Gödel’s demonstration that logic could not provide exactly the sort of foundation for mathematics or account of its formal systems that had been sought. Namely, mathematics could not be reduced to a provably complete and consistent logical theory, but logic has still remained closely allied with mathematical foundations and principles.)

Traditionally, logic had set itself the task of understanding valid arguments of all sorts, not just mathematical ones. It had developed the concepts and operations needed for describing concepts, propositions, and arguments—especially in terms of “logical form”—insofar as such tools could conceivably affect the assessment of any argument’s quality or ideal persuasiveness. It is this general ideal that many logicians have developed and endorsed, and that some, such as Hegel, have rejected as impossible or useless. For the first decades of the 20th century, logic threatened to become exclusively preoccupied with a new and historically somewhat foreign role of serving in the analysis of arguments in only one field of study, mathematics. The philosophical-linguistic task of developing tools for analyzing statements and arguments that can be expressed in some natural language about some field of inquiry, or even for analyzing propositions as they are actually (and perhaps necessarily) thought or conceived by human beings, was almost completely lost. There were scattered efforts to eliminate this gap by reducing basic principles in all disciplines—including physics, biology, and even music—to axioms, particularly axioms in set theory or first-order logic. But these attempts, beyond showing that it could be done, did not seem especially enlightening. Thus, such efforts, at their zenith in the 1950s and ’60s, had all but disappeared in the ’70s: one did not better and more usefully understand an atom or a plant by being told it was a certain kind of set.

LOGIC NARROWLY CONSTRUED

Although set theory and the type theory of Russell and Whitehead were considered to be “logic” for the purposes of the logicist program, a narrower sense of logic reemerged in the mid-20th century as what is usually called the “underlying logic” of these systems: whatever concerns only rules for propositional connectives, quantifiers, and nonspecific terms for individuals and predicates. (An interesting issue is whether the privileged relation of identity, typically denoted by the symbol “=,” is a part of logic: most researchers have assumed that it is.) In the early 20th century and especially after Tarski’s work in the 1920s and ’30s, a formal logical system was regarded as being composed of three parts, all of which could be rigorously described. First, there was the notation: the rules of formation for terms and for well-formed formulas (wffs) in the logical system. This theory of notation itself became subject to exacting treatment in the concatenation theory, or theory of strings, of Tarski, and in the work of the American Alonzo Church. Previously, notation was often a haphazard affair in which it was unclear what could be formulated or asserted in a logical theory and whether expressions were finite or were schemata standing

Relationship between set theory and logic

The theory of notation

for infinitely long wffs. Issues that arose out of notational questions include definability of one wff by another (addressed in Beth's and Craig's theorems, and in other results), creativity, and replaceability, as well as the expressive power and complexity of different logical languages.

The second part of a logical system consisted of the axioms, rules of inference, or other ways of identifying what counts as a theorem. This is what is usually meant by the logical "theory" proper: a (typically recursive) description of the theorems of the theory, including axioms and every wff derivable from axioms by admitted rules. Although the axiomatic method of characterizing such theories with axioms or postulates or both and a small number of rules of inference had a very old history (going back to Euclid or further), two new methods arose in the 1930s and '40s. First, in 1934, there was the German mathematician Gerhard Gentzen's method of succinct *Sequenzen* (rules of consequents), which were especially useful for deriving metalogical decidability results. This method originated with Paul Hertz in 1932, and a related method was described by Stanisław Jaśkowski in 1934. Next to appear was the similarly axiomless method of "natural deduction," which used only rules of inference; it originated in a suggestion by Russell in 1925 but was developed by Quine and the American logicians Frederick Fitch and George David Wharton Berry. The natural deduction technique is widely used in the teaching of logic, although it makes the demonstration of metalogical results somewhat difficult, partly because historically these arose in axiomatic and consequent formulations.

A formal description of a language, together with a specification of a theory's theorems (derivable propositions), are often called the "syntax" of the theory. (This is somewhat misleading when one compares the practice in linguistics, which would limit syntax to the narrower issue of grammaticality.) The term "calculus" is sometimes chosen to emphasize the purely syntactic, uninterpreted nature of a formal theory.

Finally, the third component of a logical system was the semantics for such a theory and language: a declaration of what the terms of a theory refer to, and how the basic operations and connectives are to be interpreted in a domain of discourse, including truth conditions for wffs in this domain. A specification of a domain of objects (De Morgan's "universe of discourse"), and of rules for interpreting the symbols of a logical language in this domain such that all the theorems of the logical theory are true is then said to be a "model" of the theory (or sometimes, less carefully, an "interpretation" of the theory).

The notion of a rigorous logical theory, in the sense of a specification, often axiomatic, of theorems of a theory, was fairly well understood by Euclid, Aristotle, and others in ancient times. With the crises in geometry of the 19th century, the need developed for very careful presentations of these theories. Hilbert's work, as well that of a group of American mathematicians that included Edward Vermilye Huntington, Oswald Veblen, and Benjamin Abram Bernstein (the postulate theorists, working shortly after 1900), reestablished this tradition with even higher standards. Frege and, in his footsteps, Russell and Whitehead, had separate claims to emphasizing standards of precision and care in the statement of logical theories. Cantor, Zermelo, and most other early set theorists did not often state the content of their axioms and theorems in symbolic form, or restrict themselves to certain symbols. Zermelo, in fact, did not often use the formal language for quantifiers and binding variables that was then available; instead, he used ordinary expressions such as "for any," "all," or "there exists." Through the 1920s, logical axioms and rules of inference were typically not all explicitly and precisely stated, especially various principles of substitution that mimicked widely understood algebraic practices.

What is known as formal semantics, or model theory, has a more complicated history than does logical syntax; indeed, one could say that the history of the emergence of semantic conceptions of logic in the late 19th and early 20th centuries is poorly understood even today. Certainly, Frege's notion that propositions refer to (German: *bedeuten*) "The True" or "The False"—and this

for complex propositions as a function of the truth values of simple propositions—counts as semantics. Earlier medieval theories of supposition incorporated useful semantic observations. So, too, do Boolean techniques of letters taking or referring to the values 1 and 0 that are seen from Boole through Peirce and Schröder. Both Peirce and Schröder occasionally gave brief demonstrations of the independence of certain logical postulates using models in which some postulates were true, but not others. (The first explicit use of such techniques seems to have arisen earlier in the 19th century, and in geometry.) The first clear and significant general result in model theory is usually accepted to be a result discovered by Löwenheim in 1915 and strengthened in work by Skolem from the 1920s. This is the Löwenheim-Skolem theorem, which states that a theory that has a model at all has a countable model. That is to say, if there exists some model of a theory (*i.e.*, an application of it to some domain of objects), then there is sure to be one with a domain no larger than the natural numbers (see above *Metalogic: Discoveries about logical calculi*). Although Löwenheim and Skolem understood their results perfectly well, they did not explicitly use the modern language of "theories" being true in "models." The Löwenheim-Skolem theorem is in some ways a shocking result, since it implies that any consistent formal theory of anything—no matter how hard it tries to address the phenomena unique to a field such as biology, physics, or even sets—can just as well be understood from its formalisms alone as being about natural numbers.

The second major result in formal semantics, Gödel's completeness theorem of 1930, required even for its description, let alone its proof, more careful development of precise concepts about logical systems—metalogical concepts—than existed in earlier decades. One question for all logicians since Boole, and certainly since Frege, had been: Was the theory consistent? In its purely syntactic analysis, this amounts to the question: Was a contradictory sentence (of the form $A \ \& \ \sim A$) a theorem? In its semantic analysis, it is equivalent to the question: Does the theory have a model at all? For a logical theory, consistency means that a contradictory theorem cannot be derived in the theory. But since logic was intended to be a theory of necessarily true statements, the goal was stronger: a theory is Post-consistent (named for the Polish-American logician Emil Post) if every theorem is valid—that is, if no theorem is a contradictory or a contingent statement. (In nonclassical logical systems, one may define many other interestingly distinct notions of consistency; these notions were not distinguished until the 1930s.) Consistency was quickly acknowledged as a desired feature of formal systems: it was widely and correctly assumed that various earlier theories of propositional and first-order logic were consistent. Zermelo was, as has been observed, concerned with demonstrating that ZF was consistent; Hilbert had even observed that there was no proof that the Peano postulates were consistent. These questions received an answer that was not what was hoped for in a later result of Gödel (discussed below). A clear proof of the consistency of propositional logic was first given by Post in 1921. Its tardiness in the history of symbolic logic is a commentary not so much on the difficulty of the problem as it is on the slow emergence of the semantic and syntactic notions necessary to characterize consistency precisely. The first clear proof of the consistency of the first-order predicate logic is found in the work of Hilbert and Wilhelm Ackermann from 1928. Here the problem was not only the precise awareness of consistency as a property of formal theories but also of a rigorous statement of first-order predicate logic as a formal theory.

In 1928 Hilbert and Ackermann also posed the question of whether a logical system, and, in particular, first-order predicate logic, was (as it is now expressed) "complete." This is the question of whether every valid proposition—that is, every proposition that is true in all intended models—is provable in the theory. In other words, does the formal theory describe all the noncontingent truths of a subject matter? Although some sort of completeness had clearly been a guiding principle of formal logical theories dating back to Boole, and even to Aristotle (and

The question of consistency

The question of completeness

to Euclid in geometry)—otherwise they would not have sought numerous axioms or postulates, risking nonindependence and even inconsistency—earlier writers seemed to have lacked the semantic terminology to specify what their theory was about and wherein “aboutness” consists. Specifically, they lacked a precise notion of a proposition being “valid,”—that is, “true in all (intended) models”—and hence lacked a way of precisely characterizing completeness. Even the language of Hilbert and Ackermann from 1928 is not perfectly clear by modern standards.

Gödel proved the completeness of first-order predicate logic in his doctoral dissertation of 1930; Post had shown the completeness of propositional logic in 1921. In many ways, however, explicit consideration of issues in semantics, along with the development of many of the concepts now widely used in formal semantics and model theory (including the term metalanguage), first appeared in a paper by Alfred Tarski, “The Concept of Truth in Formalized Languages,” published in Polish in 1933; it became widely known through a German translation of 1936. Although the theory of truth Tarski advocated has had a complex and debated legacy (see the article LINGUISTICS: *Philosophical views on meaning*), there is little doubt that the concepts there (and in later papers from the 1930s) developed for discussing what it is for a sentence to be “true in” a model marked the beginning of model theory in its modern phase. Although the outlines of how to model propositional logic had been clear to the Booleans and to Frege, one of Tarski’s most important contributions was an application of his general theory of semantics in a proposal for the semantics of the first-order predicate logic (now termed the set-theoretic, or Tarskian, interpretation).

Tarski’s techniques and language for precisely discussing semantic concepts, as well as properties of formal systems described using his concepts—such as consistency, completeness, and independence—rapidly and almost imperceptibly entered the literature in the late 1930s and after. This influence accelerated with the publication of many of his works in German and then in English, and with his move to the United States in 1939.

Gödel’s first incompleteness theorem, from 1931, stands as a major turning point of 20th-century logic. It states that no finitely axiomatizable theory sufficient to derive the Peano postulates is both consistent and complete. (How Gödel proved this fascinating result is discussed more extensively in the article MATHEMATICS, THE FOUNDATIONS OF; see also above *Metalogic: Discoveries about formal mathematical systems*.) In other words, if we try to build a theory sufficient for a foundation for mathematics, stating the axioms and rules of inference so that we have stipulated precisely what is and what is not an axiom (as opposed to open-ended axiom schemata), then the resulting theory will either (1) not be sufficient for mathematics (*i.e.*, not allow the derivation of the Peano postulates for number theory) or (2) not be complete (*i.e.*, there will be some valid proposition that is not derivable in the theory) or (3) be inconsistent. (Gödel actually distinguished between consistency and a stronger feature, ω -[omega]-consistency.) A corollary of this result is that, if a theory is finitely axiomatizable, consistent, and sufficient to derive the Peano postulates, then that theory cannot be used as a metalanguage to show its own consistency; that is, a finitely axiomatized set theory cannot be used to show the consistency of finitely axiomatized set theory, if set theory is consistent. This is often called Gödel’s second incompleteness theorem.

These results were widely interpreted as a blow to both the logicist and formalist programs. Logicists seemed to have taken as their goal the construction of rigorously described theories that were sufficient for deriving mathematics and also consistent and complete. Gödel showed that, if this was their goal, they would necessarily fail. It was also a blow to the longer-standing axiomatic, or formalist, program, since it seemed to show that precise axiomatic descriptions of valuable domains like mathematics would also necessarily fail. Gödel himself eventually interpreted the result as showing that there exist entities with well-defined properties, namely numbers, that are beyond our

ability to describe precisely with standard logical tools. This is one source of his inclination toward what is usually called mathematical Platonism.

One reply of the logicists could have been to abandon as ideal the first-order, finitely axiomatized theories, such as first-order predicate logic, the system of Russell and Whitehead, and NBG, and instead to accept theories that were less rigorously described. First-order theories allow explicit reference to, and quantification over, individuals, such as numbers or sets, but not quantification over (and hence rules for manipulating) properties of these individuals. For example, one possible logicist reply is to note that the Peano postulates themselves seem acceptable. It is true that Gödel’s result implies that we cannot prove (as Hilbert hoped in his second problem) that these postulates are consistent; furthermore, the fifth postulate is a schema or second-order formulation, rather than being strictly in the finitely axiomatizable first-order language that was once preferred. This reply, however, clashes with another desired feature of a formal theory, namely, decidability: that there exists a finite mechanical procedure for determining whether a proposition is, or is not, a theorem of the theory. This property took on added interest after World War II with the advent of electronic computers, since modern computers can actually apply algorithms to determine whether a given proposition is, or is not, a theorem, whereas some algorithms had only been shown theoretically to exist. (See the article COMPUTER SCIENCE: *Theory of computation*.) The decidability of propositional logic, through the use of truth tables, was known to Frege and Peirce; a proof of its decidability is attributable to Jan Łukasiewicz and Emil Post independently in 1921. Löwenheim showed in 1915 that first-order predicate logic with only single-place predicates was decidable and that the full theory was decidable if the first-order predicate calculus with only two-place predicates was decidable; further developments were made by Skolem, Heinrich Behmann, Jacques Herbrand, and Quine. Herbrand showed the existence of an algorithm which, if a theorem of the first-order predicate logic is valid, will determine it to be so; the difficulty, then, was in designing an algorithm that in a finite amount of time would determine that propositions were invalid. As early as the 1880s, Peirce seemed to be aware that the propositional logic was decidable but that the full first-order predicate logic with relations was undecidable. The proof that first-order predicate logic (in any general formulation) was undecidable was first shown definitively by Alan Turing and Alonzo Church independently in 1936. Together with Gödel’s (second) incompleteness theorem and the earlier Löwenheim-Skolem theorem, the Church-Turing theorem of the undecidability of the first-order predicate logic is one of the most important, even if “negative,” results of 20th-century logic.

By the 1930s almost all work in the foundations of mathematics and in symbolic logic was being done in a standard first-order predicate logic, often extended with axioms or axiom schemata of set- or type-theory. This underlying logic consisted of a theory of “classical” truth functional connectives, such as “and,” “not,” and “if . . . then,” and first-order quantification permitting propositions that “all” and “at least one” individual satisfy a certain formula. Only gradually in the 1920s and ’30s did a conception of a “first-order” logic, and of alternatives, arise—and then without a name.

Certainly with Hilbert and Ackermann’s *Grundzüge der Theoretischen Logik* (1928; “Basic Elements of Theoretical Logic”), and Hilbert’s and Paul Bernays’ minor corrections to this work in the 1930s, a rigorous theory of first-order predicate logic achieved its mature state. Even Hilbert and his coworkers, however, sometimes deviated from previous and subsequent treatments of quantification, preferring to base their theory on a single term-forming operator, ϵ , which was to be interpreted as extracting an arbitrary individual satisfying a given predicate. In the 1920s and ’30s considerable energy went into formulating various alternative but equivalent axiom systems for classical propositional and first-order logic and demonstrating that these axioms were independent. Some of these efforts were concentrated on the “implicational” (if . . . then)

Decidability

The incompleteness theorem

fragment of propositional logic. Others sought reductions of truth-functional connectives to a short list of primitive connectives, especially to the single Sheffer or, in modern terminology, NAND function, named for the American logician Henry M. Sheffer. Peirce had been aware in the 1880s that single connectives based either on not-both or on neither-nor sufficed for the expression of all truth-functional connectives. Alternative formulations of classical propositional logic reached their apex in J.G.P. Nicod's, Mordchaj Wajsberg's, and Łukasiewicz's different single-axiom formulations of 1917, 1929, and 1932. A basic underlying classical propositional logic and a first-order quantificational theory had become widely accepted by 1928, and different systems varied primarily in provably equivalent, notational aspects.

Intuitionistic logic

A notable exception to this orthodoxy was intuitionistic logic. Arising from observations by the Dutch mathematicians Arend Heyting and L.E.J. Brouwer concerning the results of indirect proof in traditional mathematics and distantly inspired by Kant's views on constructions in mathematics (and less distantly by views of French mathematicians Henri Poincaré and Émile Borel at the turn of the century), these theorists proposed that a proof in mathematics should be accepted only if it constructed the mathematical entity it talked about, and not if it merely showed that the entity "could" be constructed or that supposing its nonexistence would result in contradiction. This view is called intuitionism or sometimes constructivism, because of the weight it places on mental apprehension through construction of purported mathematical entities. (A still more severe form of constructivism is strict finitism, in which one rejects infinite sets; for further discussion, see MATHEMATICS, THE FOUNDATIONS OF: *Intuitionism*.)

The central focus of Brouwer's logical critique was directed at the principle of the excluded middle—which states that, for any proposition p , " p or not- p " is a theorem of logic—and at what one could typically infer with it. So, if not- p can be shown to be false, then in classical, but not intuitionistic, propositional logic, p is thereby proven. Intuitionistic propositional logic was formulated in 1930 by Heyting; the independence of Heyting's axioms was shown in 1939 by J.C.C. McKinsey. The primary difference between classical and intuitionistic propositional logics is concentrated in axioms and rules involving negation. Heyting in fact used the symbol \neg for intuitionistic negation, to distinguish it from the symbol \sim of classical logic.

The intuitionistic first-order predicate logic, aside from the differing propositional logic on which it is based, differs from classical first-order predicate logic only in small respects. A number of results concerning Heyting's system, as well as stronger and weaker versions of the intuitionistic propositional theory, were produced in the 1930s and '40s by the Russian theorist Andrei Nikolayevich Kolmogorov and by Mordchaj Wajsberg, Gentzen, McKinsey, Tarski, and others. Few practicing mathematicians have followed the intuitionistic doctrine of constructivism, but the theory has exerted attraction for and elicited respect from many researchers in the foundations and philosophy of mathematics. (One oddity is that metalogical results for intuitionistic logics have nearly always been shown using the theory of classical logic.)

The other major competitor to first-order predicate logic based on a classical propositional logic arose with the renewed interest in Frege's theory of properties begun by Alonzo Church in the late 1930s. The first result was a logical theory called the λ calculus, which allowed one by the application of a λ operator to a formula precisely to characterize or "extract" that property. Later developments included his investigation of formal Fregean theories ("A Logic of Sense and Denotation") that allow quantification over properties and incorporate Frege's semantic views in distinguishing between the highly individuated "sense" of an expression and its denotation (extension, or referent). These two developments laid the basis for formal theories of second- and higher-order logical theories that permit quantification over properties and other non-individuals, and for intensional logics. While Boolean and most first-order theories, including type and set theories, had dealt

Frege's theory of properties

with individuals and collections of these (collective extensions), intensional logics allow one to develop theories of properties that have the same extension but differ in intension—such as "polygons with three sides" and "polygons with three angles" or Frege's example of the morning and the evening star (*i.e.*, Venus).

Intension had often been equated with how a property is thought (its associations for the conceiver), while Frege, Church, and a number of philosophers and philosophers of language equated it with abstract, formally described entities that constitute the "meaning" or "sense" of expressions. Second-order theories and intensional logical systems have been extensively developed, and the metalogical features have been well explored. For a system like that described in Church's "A Formulation of the Logic of Sense and Denotation" (1946), consistency was shown by Gentzen in 1936, and for many similar systems it was less rigorously demonstrated by Herbrand in 1930. Weak completeness was demonstrated by Leon Henkin in 1947, although what counts as the intended interpretation and domain of such a powerful theory is problematic; strong completeness has yet to be shown, and, since it embraces first-order predicate logic, it is not decidable by a corollary of Church's own theorem. Debates about whether second-order logic is philosophically acceptable, technically usable, or even should count as "logic" in comparison with first-order theories have raged since its resurrection in the 1940s and '50s.

NONMATHEMATICAL FORMAL LOGIC

Early 20th-century formal logic was almost entirely fixated upon the project of exploring the foundations of mathematics. Furthering or exploring the logicist program and the related formalist programme of Hilbert and linking mathematics with pure logic or with rigorous formal theories had been the original motivation for many developments. The Löwenheim-Skolem theorem might have seemed also to have given a reason for this mathematical, and specifically numerical, fixation, since there is a sense in which any consistent (first-order) formal theory is always about numbers. These developments reached their height in the 1930s with the finite axiomatizations of NBG and with the formulations of the first-order predicate logic of Hilbert, Ackermann, and Gentzen. Major metalogical results for the underlying first-order predicate logic were completed in 1936 with the Church-Turing theorem. After first-order logic had been rigorously described in the 1930s and had become well understood and after set theory coalesced into ZF (with the exception of the then outstanding independence results), a period of reflection set in. There were now increasing doubts about the ability of the logicist and formalist program to connect mathematics and logic. The intuitionist critiques became well known, if not always accepted. A number of authors suggested approaching logic with entirely different formalisms—without quantifiers, for example. These included the American mathematician Haskell Curry and the category theorists, as well as algebraists who urged a return to algebraic—though not always Boolean—methods: the latter included Tarski and Paul Halmos. There were doubts about the exact form or notation and the general approach of the first-order, set-theoretic enterprise. As with many large-scale completed projects—and this project, moreover, had been accompanied by considerable disappointment, owing to the negative results of Gödel, Church, and Turing—there was also a search for new logical terrain to explore.

Łukasiewicz had, as early as 1923, begun exploring the logical theories of Aristotle and the Stoics and formalizing them as modern logical systems; this work culminated in his 1951 and 1957 editions of *Aristotle's Syllogistic* and in further work on Aristotle's logic by John Corcoran and T.J. Smiley. Benson Mates' careful study of Stoic logic similarly served to renew interest in older logics. These theories had no obvious bearing on the foundations of mathematics, but they were of interest as formal theories in their own right—and perhaps as theories of ideal reasoning, of abstract conceptual entities, or as theories of the referents of terms in natural language. Similarly, not all of Church's work on Fregean theories of properties and

Increasing doubts about the logicist program

intensions had obvious utility for constructing the simplest possible foundation for mathematics with the fewest arguable postulates, but his work was also motivated by more general theoretical features in the theory of properties and of language—especially by the richness of natural languages. These might be termed nonmathematical influences in the development of 20th-century logic. Another challenge to “classical” propositional logic—specifically to the standard interpretation of propositional logics—has been posed by many-valued logics. Propositions can be regarded as taking more than (or other than) the traditional “values” of true or false. Such possibilities had been speculated about by Peirce and Schröder (and even by medieval logicians) and were used in the 1920s and '30s by Carnap, Łukasiewicz, and others to derive independence results for various propositional calculi. In the 1940s and after, formal theories for many-valued (including infinitely valued, probabilistic-like) logics have been taken increasingly seriously—albeit for nonmathematical purposes.

Many nonmathematical goals and considerations arose from philosophy (especially from metaphysics but also from epistemology and even ethics), from the study of the history of logic and mathematics, from quantum mechanics (quantum logic), and from the philosophy of language, as well as, more recently, from cognitive psychology (starting with Jean Piaget's interest in syllogistic logics). This work has rekindled interest in logic for purposes other than giving or exploring the foundations of mathematics. Foremost among the nonmathematical interests was the development of modal logic beginning with C.I. Lewis' theories of 1932 and, specifically, a study of the alethic modal operators of necessity, possibility, contingency, and impossibility. Viable semantic accounts for modal systems in terms of Leibnizian “possible worlds” were developed by Saul Kripke, David Lewis, and others in the 1960s and '70s and led to greatly intensified research. Tense logics and logics of knowledge, causation, and ethical or legal obligation also moved rapidly forward, together with specialized logics for analyzing the “if . . . then” conditional in ordinary language (first due to C.I. Lewis as a theory of entailment, then elaborately developed by Alan Ross Anderson, Nuel Belnap, Jr., and their students as relevance logic).

From the turn of the century through the mid-1930s and with the almost singular exception of Russell and Whitehead's *Principia Mathematica*, logic was dominated by mathematicians from the German-speaking world. The work of Frege, Dedekind, and Cantor at the end of the 19th century, even if little recognized at the time, as well as the more widely recognized work of Hilbert and Zermelo, had given German mathematical logic a strong boost into the 20th century. Widespread institutional interest in the new mathematical logic in the early decades of the 20th century seemed to have been far weaker in the United States, France, and, rather surprisingly, in the United Kingdom and Italy. In the 1920s and into the early 1930s, Poland developed an especially strong logical tradition, and Polish logicians made a number of major contributions, writing in both Polish and German; in the 1920s and '30s Polish logicians posed the only exceptions to almost absolute German logical hegemony.

By the late 1930s, both the political and the logical situations had shifted dramatically. American logic, as represented by younger figures such as Church, McKinsey, and Quine, made a number of important contributions to logic in the late 1930s; the young Alan Turing in England contributed to logic and to the infant field of the theory of computation. France's place dwindled prematurely with the untimely death of Jacques Herbrand. The Moravian-Austrian Gödel fled to the United States as the political situation in central Europe worsened, as did Tarski and Carnap. Set theory and some set theorists fell under the pall of anti-Semitism, as did other logical theories, together with the theory of relativity and several philosophical orientations. Communication between scholars in Germany, both with each other and with the increasing number of researchers outside the country, was hindered in the late 1930s and '40s. With the death of Hilbert in 1943, interest in logic and in the foundations

of mathematics at the University of Göttingen—interests that had flourished there since the time of the German mathematician Bernhard Riemann—declined. With the flight of promising students and the lack of political stability and academic support, German logic became critically weakened. Heinrich Scholz, primarily a historian of logic, but also one of the few figures in Germany of the time in a philosophy, rather than a mathematics, department, attempted to rally German logic with the establishment of the Ernst Schröder Prize in mathematical logic. Its winner in 1941 was J.C.C. McKinsey, an American.

Especially because of its often predominating mathematical orientation (and this in several respects), the influence and place of 20th-century logic in all intellectual activity has changed dramatically. On the one hand, it has regained the respectability as an academic discipline through its affiliation with rigorous mathematics—the “queen of the sciences”—that it had lost in the Renaissance. On the other hand, the number of people who could profitably study modern symbolic logic and understand its more impressive achievements has dwindled as its techniques have become more austere and distant from ordinary language and have also required more and more background simply to understand. Consequently, one could say of Gödel's incompleteness theorem (for example) what Einstein once said about the theory of special relativity: that there have at times been only a handful of people who understand it. Few general university programs required an understanding of symbolic logic in the way they had once required an understanding of the rudiments of Aristotelian syllogistic or even of Venn's version of Boolean logic. Twentieth-century symbolic logic has also reexperienced its traditional problem of finding a place in modern universities.

In the early decades of the 20th century, the study of logic and the foundations of mathematics (metamathematics) acquired considerable prestige within mathematics departments, especially owing to the influence of Hilbert and the Göttingen school. In the 1920s and '30s, existing on the borderline between philosophy, mathematics, and the burgeoning interest in the philosophy of science, logic also achieved additional legitimacy through the work and participation of Wittgenstein (and Russell's philosophy of logical atomism), Carnap, and others in the Vienna and Berlin schools of scientific philosophy. (Gödel sometimes attended sessions of the Vienna Circle.)

The usefulness of logic in philosophy reached a critical point, however, with Gödel's incompleteness theorem and then with Church's logical critique of one version of the principle of verification. Roughly since the death of Hilbert, logicians and mathematical foundationalists have often been accepted less readily in mathematics departments, and after solutions to the major problems in metalogic were achieved, many practicing mathematicians in the Western Hemisphere have increasingly regarded logic and foundational work as mere tinkering. (This is less true in Russia, other former Soviet republics, and Poland, where logic has survived as a major mathematical subject.) Philosophy departments in the English-speaking world have often proved to be more stable homes for symbolic logicians, especially as they increasingly addressed issues in formal philosophy that are not necessarily issues in the foundations of mathematics, such as theories of properties and the development of nonstandard philosophical logics. (R.R.Di.)

BIBLIOGRAPHY

General works. The best starting point for exploring any of the topics in logic is D. GABBAY and F. GUENTHNER (eds.), *Handbook of Philosophical Logic*, 4 vol. (1983–89), a comprehensive reference work. See also GERALD J. MASSEY, *Understanding Symbolic Logic* (1970), an introductory text; and ROBERT E. BUTTS and JAAKKO HINTIKKA, *Logic, Foundations of Mathematics, and Computability Theory* (1977), a collection of conference papers.

Formal logic. MICHAEL DUMMETT, *Elements of Intuitionism* (1977), offers a clear presentation of the philosophic approach that demands constructibility in logical proofs. G.E. HUGHES and M.J. CRESSWELL, *An Introduction to Modal Logic* (1968, reprinted 1989), treats operators acting on sentences in first-order logic (or predicate calculus) so that, instead of being

interpreted as statements of fact, they become necessarily or possibly true or true at all or some times in the past, or they denote obligatory or permissible actions, and so on. JON BARWISE *et al.* (eds.), *Handbook of Mathematical Logic* (1977), provides a technical survey of work in the foundations of mathematics (set theory) and in proof theory (theories with infinitely long expressions). ELLIOTT MENDELSON, *Introduction to Mathematical Logic*, 3rd ed. (1987), is the standard text; and G. KREISEL and J.L. KRIVINE, *Elements of Mathematical Logic: Model Theory* (1967, reprinted 1971; originally published in French, 1967), covers all standard topics at an advanced level. A.S. TROELSTRA, *Choice Sequences: A Chapter of Intuitionistic Mathematics* (1977), offers an advanced analysis of the philosophical position regarding what are legitimate proofs and logical truths; and A.S. TROELSTRA and D. VAN DALEN, *Constructivism in Mathematics*, 2 vol. (1988), applies intuitionistic strictures to the problem of the foundations of mathematics.

Metalogic. JON BARWISE and S. FEFERMAN (eds.), *Model-Theoretic Logics* (1985), emphasizes semantics of models. J.L. BELL and A.B. SLOMSON, *Models and Ultraproducts: An Introduction*, 3rd rev. ed. (1974), explores technical semantics. RICHARD MONTAGUE, *Formal Philosophy: Selected Papers of Richard Montague*, ed. by RICHMOND H. THOMASON (1974), uses modern logic to deal with the semantics of natural languages. MARTIN DAVIS, *Computability & Unsolvability* (1958, reprinted with a new preface and appendix, 1982), is an early classic on important work arising from Gödel's theorem, and the same author's *The Undecidable: Basic Papers on Undecidable Propositions, Unsolvability Problems, and Computable Functions* (1965), is a collection of seminal papers on issues of computability. ROLF HERKEN (ed.), *The Universal Turing Machine: A Half-Century Survey* (1988), takes a look at where Gödel's theorem on undecidable sentences has led researchers. HANS HERMES, *Enumerability, Decidability, Computability*, 2nd rev. ed. (1969, originally published in German, 1961), offers an excellent mathematical introduction to the theory of computability and Turing machines. A classic treatment of computability is presented in HARTLEY ROGERS, JR., *Theory of Recursive Functions and Effective Computability* (1967, reissued 1987). M.E. SZABO, *Algebra of Proofs* (1978), is an advanced treatment of syntactical proof theory. P.T. JOHNSTONE, *Topos Theory* (1977), explores the theory of structures that can serve as interpretations of various theories stated in predicate calculus. H.J. KEISLER, "Logic with the Quantifier 'There Exist Uncountably Many,'" *Annals of Mathematical Logic* 1:1-93 (January 1970), reports on a seminal investigation that opened the way for Barwise (1977), cited earlier, and CAROL RUTH KARP, *Language with Expressions of Infinite Length* (1964), which expands the syntax of the language of predicate calculus so that expressions of infinite length can be constructed. C.C. CHANG and H.J. KEISLER, *Model Theory*, 3rd rev. ed. (1990), is the single most important text on semantics. F.W. LAWVERE, C. MAURER, and G.C. WRAITH (eds.), *Model Theory and Topoi* (1975), is an advanced, mathematically sophisticated treatment of the semantics of theories expressed in predicate calculus with identity. MICHAEL MAKKAJ and GONZALO REYES, *First Order Categorical Logic: Model-Theoretical Methods in the Theory of Topoi and Related Categories* (1977), analyzes the semantics of theories expressed in predicate calculus. SAHARON SHELAH, "Stability, the F.C.P., and Superstability: Model-Theoretic Properties of Formulas in First Order Theory," *Annals of Mathematical Logic* 3:271-362 (October 1971), explores advanced semantics.

Applied logic. Applications of logic in unexpected areas of philosophy are studied in EVANDRO AGAZZI (ed.), *Modern Logic—A Survey: Historical, Philosophical, and Mathematical Aspects of Modern Logic and Its Applications* (1981). WILLIAM L. HARPER, ROBERT STALNAKER, and GLENN PEARCE (eds.), *IFs: Conditionals, Belief, Decision, Chance, and Time* (1981), surveys hypothetical reasoning and inductive reasoning. On the applied logic in philosophy of language, see EDWARD L. KEENAN (ed.), *Formal Semantics of Natural Language* (1975); JOHAN VAN BENTHEM, *Language in Action: Categories, Lambdas, and Dynamic Logic* (1991), also discussing the temporal stages in the working out of computer programs, and the same author's *Essays in Logical Semantics* (1986), emphasizing grammars of natural languages. DAVID HAREL, *First-Order Dynamic Logic* (1979); and J.W. LLOYD, *Foundations of Logic Programming*, 2nd extended ed. (1987), study the logic of computer programming. Important topics in artificial intelligence, or computer

reasoning, are studied in PETER GÄRDENFORS, *Knowledge in Flux: Modeling the Dynamics of Epistemic States* (1988), including the problem of changing one's premises during the course of an argument. For more on nonmonotonic logic, see JOHN MCCARTHY, "Circumscription: A Form of Non-Monotonic Reasoning," *Artificial Intelligence* 13(1-2):27-39 (April 1980); DREW MCDERMOTT and JON DOYLE, "Non-Monotonic Logic I," *Artificial Intelligence* 13(1-2):41-72 (April 1980); DREW MCDERMOTT, "Nonmonotonic Logic II: Nonmonotonic Modal Theories," *Journal of the Association for Computing Machinery* 29(1):33-57 (January 1982); and YOAV SHOHAM, *Reasoning About Change: Time and Causation from the Standpoint of Artificial Intelligence* (1988).

History of logic. A broad survey is found in WILLIAM KNEALE and MARTHA KNEALE, *The Development of Logic* (1962, reprinted 1984), covering ancient, medieval, modern, and contemporary periods. Articles on particular authors and topics are found in *The Encyclopedia of Philosophy*, ed. by PAUL EDWARDS, 8 vol. (1967); and *New Catholic Encyclopedia*, 18 vol. (1967-89). I.M. BOCHENSKI, *Ancient Formal Logic* (1951, reprinted 1968), is an overview of early Greek developments. On Aristotle, see JAN LUKASIEWICZ, *Aristotle's Syllogistic from the Standpoint of Modern Formal Logic*, 2nd ed., enlarged (1957, reprinted 1987); GÜNTHER PATZIG, *Aristotle's Theory of the Syllogism* (1968; originally published in German, 2nd ed., 1959); OTTO A. BIRD, *Syllogistic and Its Extensions* (1964); and STORRS MCCALL, *Aristotle's Modal Syllogisms* (1963). I.M. BOCHENSKI, *La Logique de Théophraste* (1947, reprinted 1987), is the definitive study of Theophrastus' logic. On Stoic logic, see BENSON MATES, *Stoic Logic* (1953, reprinted 1973); and MICHAEL FREDE, *Die stoische Logik* (1974).

Detailed treatment of medieval logic is found in NORMAN KRETZMANN, ANTHONY KENNY, and JAN PINBORG (eds.), *The Cambridge History of Later Medieval Philosophy: From the Rediscovery of Aristotle to the Disintegration of Scholasticism, 1100-1600* (1982); and translations of important texts of the period are presented in NORMAN KRETZMANN and ELEONORE STUMP (eds.), *Logic and the Philosophy of Language* (1988). For Boethius, see MARGARET GIBSON (ed.), *Boethius, His Life, Thought, and Influence* (1981); and for Arabic logic, NICHOLAS RESCHER, *The Development of Arabic Logic* (1964). L.M. DE RIJK, *Logica Modernorum: A Contribution to the History of Early Terminist Logic*, 2 vol. in 3 (1962-1967), is a classic study of 12th- and early 13th-century logic, with full texts of many important works. NORMAN KRETZMANN (ed.), *Meaning and Inference in Medieval Philosophy* (1988), is a collection of topical studies.

A broad survey of modern logic is found in WILHELM RISSE, *Die Logik der Neuzeit*, 2 vol. (1964-70). See also ROBERT ADAMSON, *A Short History of Logic* (1911, reprinted 1965); C.I. LEWIS, *A Survey of Symbolic Logic* (1918, reissued 1960); JØRGEN JØRGENSEN, *A Treatise of Formal Logic: Its Evolution and Main Branches with Its Relations to Mathematics and Philosophy*, 3 vol. (1931, reissued 1962); ALONZO CHURCH, *Introduction to Mathematical Logic* (1956); I.M. BOCHENSKI, *A History of Formal Logic*, 2nd ed. (1970; originally published in German, 1962); HEINRICH SCHOLZ, *Concise History of Logic* (1961; originally published in German, 1959); ALICE M. HILTON, *Logic, Computing Machines, and Automation* (1963); N.I. STYAZHKIN, *History of Mathematical Logic from Leibniz to Peano* (1969; originally published in Russian, 1964); CARL B. BOYER, *A History of Mathematics*, 2nd ed., rev. by UTA C. MERZBACH (1991); E.M. BARTH, *The Logic of the Articles in Traditional Philosophy: A Contribution to the Study of Conceptual Structures* (1974; originally published in Dutch, 1971); MARTIN GARDNER, *Logic Machines and Diagrams*, 2nd ed. (1982); and E.J. ASHWORTH, *Studies in Post-Medieval Semantics* (1985).

Developments in the science of logic in the 20th century are reflected mostly in periodical literature. See WARREN D. GOLDFARB, "Logic in the Twenties: The Nature of the Quantifier," *The Journal of Symbolic Logic* 44:351-368 (September 1979); R.L. VAUGHT, "Model Theory Before 1945," and C.C. CHANG, "Model Theory 1945-1971," both in LEON HENKIN *et al.* (eds.), *Proceedings of the Tarski Symposium* (1974), pp. 153-172 and 173-186, respectively; and IAN HACKING, "What is Logic?" *The Journal of Philosophy* 76:285-319 (June 1979). Other journals devoted to the subject include *History and Philosophy of Logic* (biannual); *Notre Dame Journal of Formal Logic* (quarterly); and *Modern Logic* (quarterly). (M.L.Sc./P.V.S./R.R.D.)

London

The capital city of the United Kingdom, London lies astride the River Thames 50 miles (about 80 kilometres) upstream from its estuary on the North Sea. It has a population of about 6.6 million. In satellite photographs the metropolis can be seen to sit compactly in a Green Belt of open land, with the M25 orbital motorway threaded around it at a radius of about 20 miles from the city centre. The growth of the built-up area was halted by strict town planning controls in the mid-1950s. Its physical limits more or less correspond to the administrative and statistical boundaries separating London from the "Home Counties" of Kent, Surrey, and Berkshire (in clockwise order) to the south of the river and Buckinghamshire, Hertfordshire, and Essex to the north.

If the border of the metropolis is well defined, its internal structure is immensely complicated and defies description. Indeed, London's defining characteristic is an absence of overall form. It is physically a polycentric city, with many core districts and no clear hierarchy among them. London has at least two (and sometimes many more) of everything: cities, mayors, dioceses, cathedrals, chambers of commerce, police forces, opera houses, orchestras, and universities. In every aspect it functions as a compound or confederal metropolis.

Historically, London grew from three distinct centres: the walled settlement founded by the Romans on the banks of the Thames in the 1st century AD, today known as "the Square Mile" or simply "the City"; facing it across the bridge on the lower gravels of the south bank, the suburb of Southwark; and a mile upstream, on a great southward bend of the river, the City of Westminster. The three settlements had distinct and complementary roles. London, "the City," developed as a centre of trade, commerce, and banking. Southwark, "the Borough," became known for its monasteries, hospitals, inns, fairs, pleasure houses, and the great theatres of Elizabethan London—the Rose (1587), the Swan (1595), and Shakespeare's Globe (1599). Westminster grew up around an abbey, which brought a royal palace and, in its train, the entire central apparatus of the British state—its legislature, executive, and judiciary. It also boasts spacious parks and the most fashionable districts for living and shopping—"the West End." The north-bank settlements merged into a single built-up area in the early decades of the 17th century, but they did not combine into a single enlarged municipality. The City of London was unique among Europe's capital cities in retaining its medieval boundaries. Westminster and other suburbs were left to develop their own administrative structures—a pattern replicated a hundred times over as London exploded in size, becoming the prototype of the modern metropolis.

The population of London already exceeded one million by 1800. A century later it reached 6.5 million. The city's physical expansion was not constrained either by military defenses (a highly influential factor on mainland Europe) or by the intervention of state power (so evident in the

town planning of Paris, Vienna, Rome, and other capitals of continental Europe). Though much of the land around London was owned by the aristocracy, the church, and other institutions with feudal roots, its development was the work of unfettered capitalism, driven by the housing demands of the rising middle class. Free-ranging building speculation engulfed villages and small towns over an ever-widening radius with each improvement in transport technology and purchasing power. The solidly built-up area of London measured some 5 miles from east to west in 1750, 15 in 1850, and 30 in 1950.

The evacuation and bombing during World War II were a turning point in London's history because they brought the long era of expansive suburbanization to a sudden end. It was decided by the government that the metropolis had grown too much for its own economic and social good and that its growth was a strategic risk. A Green Belt was imposed after the war, and subsequent growth was diverted beyond it. Later, London's administrative boundaries were redrawn to incorporate almost the entire physical metropolis, an area of 610 square miles (1,580 square kilometres).

The London known to international visitors is a much smaller place than that. Tourist traffic concentrates on an area defined by the main attractions, each drawing between one and seven million visitors in the course of the year: the British Museum, the National Gallery, Westminster Abbey, Madame Tussaud's waxwork collection, the Tower of London, the three great South Kensington museums (Natural History, Science, and Victoria and Albert), and the Tate Gallery. In scale, visitors' London resembles the metropolis as it was in the late 18th century, a city of perhaps 10 square miles, explorable on foot in all directions from Trafalgar Square.

Resident Londoners see the metropolis in even more localized terms. Property correspondents and estate agents like to describe London as a collection of villages, and there is some truth in their cliché. Because London had developed in a dispersed, haphazard fashion from an early stage, many of its later suburbs were able to grow around, or within reach of, some existing nucleus such as a church, coaching inn, mill, parkland, or common. Buildings of different ages and types help to define the character of residential areas as well as to relieve suburban monotony. The population in the various neighbourhoods tends to be diverse because the working of the English housing market has provided most areas, even the most exclusive, with at least some public rental housing. The chemistry of location, building stock, local amenities, and property values combines with that of a multiethnic population to give rise to a great variety of residential microcosms within the metropolis. Neighbourhood ties are strong. Wherever Londoners meet and talk, they avidly compare nuances of the districts they live in because where they live seems to count for as much as who they are.

This article is divided into the following sections:

Physical and human geography 284

The landscape 284

Site

Climate

Environment

The city layout

The people 286

The economy 287

Trade, administration, and leisure

Industry

Finance

Transportation

Administration and social conditions 289

Cultural life 293

History 294

The early period 294

Foundation and early settlement

Medieval London

Tudor London

17th-century London

Evolution of the modern city 297

18th-century London

Organization, innovation, and reform

Reconstruction after World War II

Bibliography 298

Physical and human geography

THE LANDSCAPE

Site. *The geologic foundation.* The landscape of south-eastern England is shaped by an undulating bed of thick white chalk, consisting of a pure limestone speckled with flint nodules in the upper beds. Under the chalk are an incomplete layer of Upper Greensand (a Cretaceous rock; 66.4 to 144 million years old) and a 200-foot- (60-metre-) thick waterproof layer of Gault clay. Beneath them in turn lies London's true geologic foundation, a stable platform of old hard rocks of Paleozoic age (about 245 to 540 million years old). This basement is buried nearly 1,000 feet below London, sloping away southward to depths more than 3,300 feet below the English Channel.

The London Basin is a wedge-shaped declivity bounded to the south by the chalk of North Downs, running north to south, and to the north by the chalk outcrop of the Chiltern Hills, running up in a northeasterly direction from the Goring Gap. The chalk floor of the basin carries a sequence of clays and sands of the Tertiary Period (those 1.6 to 66.4 million years old), chiefly the stiff, gray-blue London Clay, which lies up to 433 feet thick under the metropolis and supports most of its tunnels and deeper foundations. The subsoil is topped with deposits of gravel up to 33 feet deep, consisting mostly of pebbles with flint, quartz, and quartzite. There are also patchy deposits of brick-earth, a mixture of clay and sand often excavated for building materials. Lastly, modern London is built on "made ground," the deposits of centuries of continuous human occupation, which have accumulated on average between 10 and 16 feet in the oldest urban nuclei of the City and Westminster.

The valley of the Thames. The metropolis grew and spilled over a more or less symmetrical valley site defined by shallow gravel and clay ridges rising to about 450 feet on the north at Hampstead and about 380 feet at Upper Norwood 11 miles to the south. Between these broken heights to the north and south, the ground falls away in a series of graded plateaus formed by gravel terraces—some at 100–150 feet (the Boyn terraces, such as Islington, Putney, and Richmond) and a second and more extensive level, the Taplow terraces, at 50–100 feet, on which sit the City of London, the West and East Ends, and the elevated southern districts such as Peckham, Battersea, and Clapham. The lowest ground, just a few feet above high-tide level, is the extensive floodplain of the valley floor. The Thames scours the confining terraces to the north and south as it meanders toward the sea. The Romans founded the city of London where the northernmost meander undercuts the higher gravel terrace to form a steep bluff. Here, at the upper limit of tidal navigation, was an ideal location for defense and commerce alike. Most of London's subsequent growth extended from this nucleus along the better-drained terraces of the north bank. Building remained more difficult in the alluvial ground south of the river until the completion of tidal embankments in the 19th century.

To complete the picture of London's site in its natural state before building took place, one must add the tributary streams running north and south from the hills to the great river on the valley floor, many of them rising from springs in the gravel. Those in the centre of town have long since been culverted over, except where they do duty as ornamental water in parks (e.g., the Serpentine in Hyde Park). Their names survive in the topography of London: Holborn, Fleet Street, Walbrook. Away from central London are a series of larger tributaries, used variously for navigation and associated industries, drinking water collection, gravel extraction, and ornament and recreation. To the northwest the River Colne and the River Crane join the Thames at Staines and at Isleworth, respectively; to the northeast the Lea, a substantial river draining much of Hertfordshire, enters the Thames just beyond the Isle of Dogs at Blackwall; and the River Roding merges into it four miles downstream at Barking. South London has a series of much smaller rivers leading north to the mainstream: the Ravensbourne flows through Bromley, Lewisham, and Deptford, entering the

tidal Thames at Greenwich; the River Wandle rises near Croydon and flows down through Merton and Tooting to join the Thames at Wandsworth; Beverley Brook rises in Sutton and runs at the foot of Wimbledon Common and through Richmond Park and Barnes Common, emerging from a culvert at Barn Elms; the Hogsmill River flows down from the Epsom Downs to Kingston-upon-Thames; and, in the southwest corner of modern London, the River Mole drains the Surrey hills to join the Thames opposite Hampton Court.

Panorama of the city. The natural lay of the land can be appreciated from several public vantage points. Hampstead Heath offers the finest panorama over the central basin of the metropolis. But from Shooters Hill, Upper Norwood, or Alexandra Palace one has a choice of views: inward to the crowded skyline of the City and West End or out to the open expanses of the Home Counties, the Thames estuary, the South Downs, and the Weald. Such panoramas show that London, for all its immensity, resembles more closely the limited metropolises of the early 20th century than the amorphous and sprawling megapolises of today, such as Tokyo or Los Angeles. The line of the post-World War II Green Belt runs quite comfortably along the encircling hills of the London Basin—the long ridge of the downs to the south of London and, to the north, the more broken chain of heights running from Iver Heath (above Heathrow Airport) clockwise through Ruislip Common, Bushey Heath, Enfield Chase, Epping Forest, Hainault Forest, and South Weald.

Climate. Continuous records of London's weather extend back to 1659, with specific data for wind direction available since 1723 and for rainfall since 1697. The fluctuations show a cyclic pattern, with troughs of hard winters and cold springs during the 1740s, 1770s, 1809–17, 1836–45, and 1875–82 followed by a long upswing after 1919, in which London's climate became warmer, largely because of milder weather in the autumn months.

Modern London has the equable climate of South East England, with mild winters and temperate summers. The average daytime air temperature is 52° F (11° C), with 42° F (5.5° C) in January and 65° F (18° C) in July. Statistics show that the sun shines, however briefly, on five days out of six. Londoners shed their winter overcoats in April or May and begin to dress warmly again in late October. The prevailing wind is west-southwest. Because of the sheltering effect of the Chiltern Hills and North Downs, the city has slightly less rainfall than the Home Counties. In an average year one can expect 200 dry days out of 365 and a precipitation total of about 23 inches quite evenly distributed across the 12 months.

The incidence of sleet and snow is less predictable. It varies greatly from year to year around a long-run statistical average of 20 days. The snowiest winter on record was 1695, with snow falling on 70 days. When snow does fall (generally only in the first three months of the year), it rarely accumulates. Semihardy plants can winter over in London gardens, though only in the most sheltered and sunny spot will a London vine bear grapes sweet enough for wine making.

Climatic variations across the metropolis show very clearly that there is a heat island created by concentration of buildings, internal-combustion engines, and heating and air-conditioning plants. Temperatures are higher toward the centre of the city, and the air is drier. Overall, the average difference in minimum temperatures between London and the surrounding country is 3.4° F (1.9° C), but on individual nights the difference can be as much as 16.2° F (9° C). The chemical, mechanical, and thermal effects of the city also affect wind speed and precipitation. Downpours of heavy rain are liable to be more intense within London because pollution particles act as condensation nuclei for water vapour.

Environment. *Smog and air pollution.* For years London was synonymous with "smog," the word coined at the turn of the 20th century to describe the city's characteristic blend of fog and smoke. The capital's "pea-soupers" were caused by suspended pollution of smoke and sulfur dioxide from coal fires. The most severely affected area was the 19th-century residential and industrial belt of inner

Records of
London's
weather

Tributaries
of the
Thames

Clean Air Acts

London—particularly the East End, which had the highest density of factory smokestacks and domestic chimney pots and the lowest-lying land, inhibiting dispersal. As recently as the early 1960s, the smokier districts of east Inner London suffered a 30 percent reduction in winter sunshine hours. That problem was alleviated by parliamentary legislation (the Clean Air Acts of 1956 and 1968) outlawing the burning of coal, combined with the clearance of older housing and the loss of manufacturing.

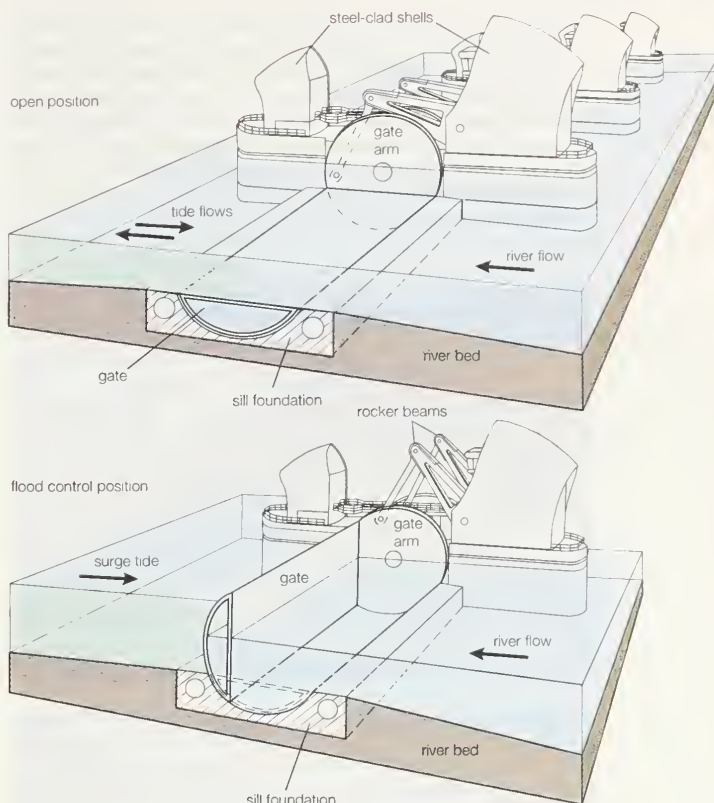
The less visible but equally toxic pollutants of carbon monoxide, nitrogen dioxide, ozone, benzines, and aldehydes continue to spoil London's air. Traffic fumes and other exhausts are liable to become trapped between the surrounding hills and below a stagnant capping mass of warm urban air at 3,000 feet, causing immediate increases in eye irritation, asthma, and bronchial complaints. But London's weather is too fickle for the development of a full-scale photochemical smog of the kind that can build up under the more stable weather conditions of cities such as Los Angeles.

Water pollution. Until the 1960s the waters of London's rivers were as polluted as its air. Deoxygenated and black with scum, they showed the effects of sewage pollution and uncontrolled industrial effluents. Then tighter environmental standards combined with the closure of factories to produce an improvement in water quality. Salmon, sea trout, roach, and flounder returned to the tidal Thames, together with shrimps, prawns, sea horses, and (at the other end of the size range) giant conger eels. Large-scale fishing of eels, a traditional Cockney delicacy, was restarted after a hiatus of 150 years. Herons, cormorants, gannets, grebes, shelducks, pochards, and terns recolonized the river.

Flood control. The greatest concern in the management of the River Thames has been the risk of flooding. Its waters are rising at the rate of 2.8 feet per century. The record floods of 1791 reached a height of 14 feet above the fixed measuring point, Ordnance Datum at London Bridge; those of 1953 rose to 17.7 feet. At high tide on a spring day, when the river is swollen with rain, it is awesome to see ships moored along the Victoria Embankment riding high above the roadway, and it is sobering to reflect on the damage that would result if the waters overtopped the walls. A serious flood would threaten 45 square miles of London's low-lying land, affecting one and a quarter million people and a quarter of a million buildings and paralyzing the capital's dense infrastructure of underground railways; sewers; gas, water, and electricity mains; telephone cabling; and service tunnels.

The flood risk results from a combination of factors. All of southeastern Britain is being tilted down into the sea (and the Hebrides tilted up) by tectonic movements resulting from the melting of Pleistocene ice sheets. London is sinking faster than the remainder of the region because modern techniques of water extraction from the chalk aquifer have dried up the underlying beds of clay. Meanwhile, sea levels in general are rising because of the effect of global warming on the polar ice caps. In addition, the tidal rhythm of the Thames has been amplified by dredging for navigation and by the embankment of its estuary marshes for cultivation.

The traditional method of protection was to build up the river walls and embankments. Long stretches were raised after the Thames Flood Act of 1879; further protective measures were taken after serious flooding in 1928, when 14 people drowned in basements in Westminster, and again after the still more serious inundations in 1953. The official inquiry into the 1953 floods recommended that "apart from erecting further walls and banks, an investigation should be made into the building of a flood barrier across the Thames." After 20 years of debate about the best design and location for a barrier, the Greater London Council settled on an unusual form of flood protection that leaves the tidal Thames intact. At Silvertown, eight miles downstream of London Bridge, it erected a line of piers from which are suspended 10 mighty steel gates and counterweights, the 4 main ones weighing 3,000 tons each. Normally laid face-downward on the bed of the river, at a time of flood risk they can be swung up by



The Thames Barrier consists of 10 movable gates separated by 9 piers. Each gate has a curved face that lies in a recessed chamber in the riverbed when the barrier is fully open. When the signal is given, the gates rotate 90° to a closed position, blocking the path of the surge tide in less than 30 minutes.

Encyclopædia Britannica, Inc.

electrohydraulic machinery to form a continuous barrier sealing London off from the sea. Downstream of the Thames Barrier, to protect against the backsurge caused by its closure, elaborate walls were built along the estuary marshes with guillotine-style floodgates at the mouths of tributary rivers.

The city layout. *Three basic patterns.* London's complicated topography can be made simple by means of three basic patterns. First, there is the wiggling line of the Thames separating north from south London. For historical reasons, most important destinations lie north of the river. The south is essentially an intricate patchwork of residential districts joined by miles of conventional through streets. It has no fast through roads.

Second, London differs from east to west. The waters of

J. Messerschmidt/Tony Stone Images



Piccadilly Circus in the West End.

the Thames and the prevailing winds flow eastward. Therefore, shipping, heavy haulage, manufacturing, and labouring districts developed downstream in the East End, while the affluent and leisured classes built their homes and pursued their pleasures in the West End. This social gradient was reinforced by the location of the royal palaces at Westminster, Kensington, Richmond, and (beyond London's boundary) Windsor. Partly in consequence, the western sector has a series of tranquil and elegant open spaces to either side of the river, from St. James's Park, by the prime

(Top) Andrew Butler, (bottom) Envision/E. V.K. Guy Photo Library



(Top) Row houses in Bedford Square, Bloomsbury. (Bottom) The Pond Gardens, Hampton Court.

minister's house at No. 10 Downing Street, through Hyde Park, Kensington Gardens, Battersea Park, Wimbledon Common, Richmond Park, the Royal Botanic Gardens at Kew, the Richmond riverbank, Hampton Court Park, and Bushey Park. Their landscapes soften the effect of noise pollution under the flight path of Heathrow Airport, on the western border. Proximity to the world's busiest international airport has itself reinforced the favoured position of western London.

The east-west divide is entrenched equally in the physical fabric of London and in the psychology of Londoners. Its significance did, however, diminish in the later years of the 20th century as the port economy and manufacturing industry declined and as white-collar work and residents began taking their place. This process became accelerated after 1981 by a generously funded government regeneration project for the 5,000 derelict acres of the London Docklands.

Overlying the north-south and east-west distinctions is a simple concentric ring pattern that reflects the historical phases of London's growth. At the centre is the area familiar to visitors, with its offices, shops, and public buildings. Next comes the suburban belt—known for statistical purposes as Inner London—developed from the late 18th century until the beginning of World War I. Terraced houses predominate, and the building scale is domestic and intimate, except where the original units were replaced by higher-density rental housing built by local councils in areas of bomb damage or postwar clear-

ance. The third concentric zone—Outer London—consists of 20th-century suburban housing, chiefly created in a short, intensive building boom in 1925–39. The most common building type is the semidetached unit, a distinctively British compromise between row housing and the freestanding homestead. The Metropolitan Green Belt forms a final concentric ring, defining the shape of the whole capital.

Population density. The metropolis has an overall population density of just above 10,500 persons to the square mile (just over 4,000 persons to the square kilometre), which is about average by British standards (Birmingham's is 9,900 and Liverpool's 11,000) but only one-fifth that of Paris. The 20 boroughs of Outer London have an average of 8,000 people to the square mile, while the 13 inner boroughs have 21,200. Yet even in Inner London the pattern of the streets and the style of the housing lacks the intense urban density of the great cities of mainland Europe. Only one residence in three was meant to be an apartment house. More than half of London's dwellings are houses with their own patch of land. The most common type is the terraced, or row, house. Monumental and institutional buildings take their place in a loose and predominantly residential urban fabric that leaves much land unbuilt even in the areas of densest development. The city's architecture is individualistic and variable, reflecting a political aversion, in this bourgeois metropolis, to the imposed order of a set piece. Only rarely are buildings used as component parts of a larger townscape composition.

THE PEOPLE

Demographic trends. From a total population of 5,600,000 in 1891, London grew by 3,000,000 to its peak magnitude at the outbreak of World War II. For the remainder of the 20th century its population shrank by approximately 2,000,000, an average loss of 40,000 people per year. Its demographic decline occurred for reasons common to all large cities of its type. Increasing leisure and holiday time, shorter working hours, and access to the automobile freed people from the ties of proximity to their place of work. Families moved out of town in search of a better quality of life. Firms moved for similar reasons, seeking more spacious and accessible sites. As the remaining population spread itself more comfortably in the dwelling stock, the three-generation household became a rarity except among ethnic minorities. Mass housing initiatives and individual "gentrification" of terraced houses tended equally to reduce population density.

The steepest fall occurred in the densest areas. Inner London boroughs lost more than a third of their population in the decades after World War II. In the 1980s the slump was eased by a fall in the rate of out-migration and a rise in the birth rates of new immigrant families. At the time of the 1991 census, one birth in three was to a mother born outside the United Kingdom. London's birth rate stood somewhat above the national average, and its mortality rate was lower. Its population was stabilizing at about 6.6 million, almost comparable in size to that of New York City, which also happens to sit in a wider urban region of approximately 20 million people.

Ethnic composition. *The historical base.* The stabilization of total population numbers masked a continuing population flux. London, like any great metropolis, acts as a nursery, perpetually taking in young and aspiring immigrants and releasing mature firms and families. But the hinterland has shifted. In the 19th century most movement into London was domestic; the majority of immigrants came from the neighbouring Home Counties, with additional long-distance streams from Wales, Ireland, and Scotland. Overseas immigrants came too, but London was less cosmopolitan than New York City or Boston. Its alien communities were small (mostly less than 1,000 people) and localized, and some were long-established. Bevis Marks, the City synagogue of the Sephardic Jews, was founded in 1656. St. Peter's Italian Church (1863) was the first Italian church ever to be built outside Italy.

Immigrants from Europe. In late 19th-century London, Italians clustered in Holborn and Finsbury, French in Soho, and Chinese near the docks in Limchouse, and there

Population decline

London's historic ring pattern

was a scattering of Germans and Scandinavians around the City. Communities of Irish (at that time still subjects of the British crown) were established in Wapping and Camden. The pogroms of the 1880s and '90s brought about 20,000 Polish and Russian Jews to settle on the eastern edge of the City at Whitechapel. A further wave of Jewish emigrants fled to London from German fascism in the 1930s, followed by a wave of refugees from central Europe in the upheavals at the end of World War II.

After the war, the Polish community sank its roots in Ealing in western London. Jewish families became suburban, concentrating especially in Edgware, Golders Green, Hendon, and Finchley to the northwest and Ilford to the northeast. The extreme orthodox did not move as far, only to the northern edge of the East End in Hackney. Sizable communities of Greek and Turkish Cypriots arrived to open shops, restaurants, and small businesses on the City fringe, rising rapidly to suburban prosperity along the radial roads northward.

London's
black
population

Immigrants from the Commonwealth. London's black population grew significantly during the economic boom years of the 1950s and '60s. It was a time of labour shortage, particularly for public service industries such as transportation (buses and underground) and hospitals. To fill the places of blue-collar workers who had been encouraged to leave London to take jobs in the new towns, employers began to recruit from the former colonies, which were now independent members of the British Commonwealth. The first wave of immigration was from the Caribbean. Black Londoners found it hard to gain access to public rental housing, and so concentrated as private tenants in lodging-house districts of North Kensington and south of the river in Brixton. Kensington's Notting Hill Carnival, at the end of August, remains the chief annual celebration of West Indian life in London. Later groups of immigrants from the Commonwealth settled in different parts of the city: Indians in Ilford, Ealing, and Hounslow; Bangladeshis in Whitechapel (where they replaced the Jews in an unusually exact immigrant succession); and Africans in Hackney, Southwark, and Lambeth.

The multiethnic metropolis. London, always a cosmopolitan city, grew steadily more polyglot and multicultural. The Commonwealth connection accounted for only part of the transformation. Despite restrictive immigration laws, the flux of refugees and asylum-seekers from many countries continued, and new communities of Vietnamese, Kurds, Somalis, Eritreans, Iraqis, Iranians, Brazilians, and Colombians sprang into being. Many of the foreigners settled into hard-to-let housing estates in the poorer parts of Inner London, particularly the crescent of inner boroughs to the east of the City. At the other end of the economic spectrum, London's position at the crossroads of the global economy brought transient populations of the international business world as well as the schools, shops, and renting agencies and services to support them. Their social geography was entirely different, spreading in an arc through the northwest and southwest suburbs. London also attracted wealthy foreigners to become property owners and seasonal residents. Thus, people from the Middle East, East Asia, and Latin America purchased real estate and internationalized neighbourhoods such as Mayfair, Park Lane, and Belgravia. Shopping streets that lead north from Hyde Park, such as Queensway and the southern end of Edgware Road, are almost entirely taken over by Arabs.

Though it is not easy to establish reliable figures on London's ethnic composition, the columns of names in the telephone books and school registers are testimony to the transformation of a population that in the middle years of the 20th century had still been chiefly British-born and Anglophone. Just under one-quarter of the resident population of modern London comes from overseas. The transformation of London into a multiethnic city has gone furthest in the western boroughs (partly because of the proximity to Heathrow), while the boroughs of Havering, Barking and Dagenham, Bexley, and Bromley form an arc of almost entirely British-born white population on the far eastern edge of London. These are also the areas least touched by the cosmopolitan restaurants, clubs, and

shops that have banished the old, insular dining habits elsewhere in the metropolis.

Residential patterns. London's social geography is never static. The city has never had ghettos or strong policies of segregation. The areas of local government are too large and the housing stock too diverse for exclusionary practices of the kind encountered in some North American cities. There is intermixture even in the areas having a high concentration of one particular group, such as those of the extreme orthodox Jews at Stamford Hill, the Sikhs at Southall, or the West Indians at Brixton. Boundaries and distributions are perpetually shifting. Minorities follow one another in the familiar sequence of arrival, consolidation, and outward and upward mobility. Jews who came to Whitechapel in the 1890s shifted eastward to the semi-detached suburb of Ilford. Cypriots who had settled along the Seven Sisters Road moved north along the old drovers' road, Green Lanes, to Tottenham and Haringay. Traces of earlier diasporas are scattered through Inner London. Most of London's 11 Welsh churches are grouped around the centre. The Welsh Congregational Church at Radnor Walk in Chelsea today serves a dispersed instead of a local congregation. Swedish, Norwegian, and Danish Lutherans drive eastward on Sunday mornings to worship in their old churches at the dockyard gates.

London's
social
geography

THE ECONOMY

Trade, administration, and leisure. London has been described above as a polycentric city. The map of Elizabethan London shows that fields and the river separated distinct centres: the City of London with its shipping, trade, and crafts; Southwark with its gardens, hospitals, and theatres; and the royal court at Westminster. The economy of modern London has evolved continuously from the three complementary elements of trade, administration, and leisure. In trade London is one of a handful of centres—along with New York City, Tokyo, and Hong Kong—where dealers in currencies, equities, commodities, and insurance operate on a global scale. In the first half of the 20th century it was also a substantial manufacturing centre. In contrast to the other great cities of Britain, London's factory closures have been compensated at least partly by the city's dynamism in financial services and the media.

As an administrative centre, London dominates the national life to an exceptional degree. The United Kingdom is constitutionally a unitary state and politically the most centralized in Europe. Scotland, Wales, and Northern Ireland, England's three national partners within the United Kingdom, have administrative identity but lack political institutions. All legislative activity is concentrated in the English capital, at Westminster. Pressure groups and lobbyists needs must follow. British local governments raise less than one-quarter of their needs in tax revenues and are heavily dependent on fiscal transfers from the centre. In British politics, all roads lead to London.

If London is a place to win influence and make money, it is also a great playground—a leisure metropolis. Historically the landed classes flocked to London each year to spend "the season" in the proximity of the court. The legacy of aristocratic consumption still survives in the gunsmiths, art dealers, tailors, and vintners of the West End, serving a modern market of London's international visitors. Each year more than 100 million nights are spent by tourists in the capital's hotels. Though its full impact is difficult to trace, the tourism industry has clearly overtaken manufacturing as a source of employment for Londoners, offering direct employment for more than 200,000 workers and perhaps as many more again through economic multiplier effects, some of them in the black market.

Industry. Shipping. For centuries, shipping was at the heart of the economy of London. The city retained its lead as the largest, busiest port in the world until World War II, with an average of 1,000 ship arrivals and departures every week. The Port of London Authority, founded in 1909, supervised seven systems of enclosed docks with a combined water area of 720 acres (291 hectares). It had 35 miles of dock quays and as many again of riverside moorings, wharfage, shipyards, and heavy industry along the

banks of the Thames from Gravesend to London Bridge.

Shipping left London quite suddenly between 1968 and 1981 for a combination of reasons, including the containerization of ocean traffic and the growing scale of bulk cargoes, poor labour relations, and competition from new private ports based in small towns around the coast. After consolidating its activities in docks at Tilbury on the Thames estuary 26 miles downstream of London, the Port of London Authority was left with a diminished share of about 8 percent of the nation's total port traffic.

Manufacturing. In addition to its importance in administration and banking, London used to be a substantial manufacturing centre. In the 18th and 19th centuries its industries were quite comparable to those of other European capitals and court cities, producing such luxury items as silks, fine furniture, gilded work, watches, musical instruments, millinery, and women's clothing. Such highly skilled trades with their own systems of apprenticeship clustered tightly around the City of London and adjacent districts. In the 20th century London became the preferred location for a new generation of electrically powered industries serving mass consumption markets. Many of the companies were American multinationals, including Heinz Company, Hoover, Ford Motor Company, and Firestone, while others had grown out of conventional craft industries. Their factories, often built in a hygienic white-tiled Art Deco style, lined the new arterial roads out of London: the Great West Road, Western Avenue, and Purley Way. In both new and old sectors, London's manufacturing base rested upon industries producing consumer goods (rather than intermediate and capital goods) such as leather products, clothing, timber and furniture, food and drink, pharmaceuticals, and specialized goods as well as products generated by printing and publishing, instrument engineering, and electrical engineering.

This manufacturing success in London presented such a striking contrast to the high levels of unemployment in the old, established industrial regions of northeastern England and Clydeside (Scotland) that the government, fearing massive expansion of the metropolis, decided to halt the city's growth. It did so by imposing a Green Belt, or stopline, to keep postwar London within strict bounds. Many growth industries, with their young and skilled workforces, were relocated to public satellite towns. Grants and incentives attracted firms into peripheral regions of high unemployment, while administrative controls discouraged the building of factories in London. At the local level, there was a tendency in the 1950s and '60s to sweep away older mixed industrial districts in urban renewal programs, while scattered industrial premises were weeded out as "non-conforming uses" under zoning powers.

Even without such public discouragement, London's industry would probably have declined in modern times because of wider shifts in the geography of manufacturing. Firms have tended to move out of large cities everywhere, drawn by access to the national highway network and by the flexibility and efficiency of low-density units on greenfield sites. The deindustrialization of London was a drawn-out and painful process. In the 1950s the decline of older craft-based manufacturing in the inner parts was concealed or compensated by continued growth of the newer industries of the interwar period. As late as 1961 half of the jobs in the London suburbs were in manufacturing. Thereafter, however, the curve of employment in manufacturing sloped remorselessly downward. A third of a million jobs were lost in the 1960s, almost half a million in the 1970s, and a further third of a million in the 1980s. Toward the end of the 20th century, Britain's greatest manufacturing city had become a "postindustrial" metropolis, with only a residual one-tenth of its workforce in manufacturing.

The main surviving concentrations of industry in London were along transport corridors. The first and still the foremost of these are the River Thames and (to a lesser extent) its tributaries, especially for industries linked to sea-carried bulk cargoes such as petrochemicals, sugar, grains, and timber. Industrial plants continued to dominate the riverside landscape downstream of Greenwich, but upstream they were almost entirely replaced by res-

idential apartments and office blocks. Other significant manufacturing districts were on the arterial roads leading out of London and around the North Circular Road that rings London five miles from the centre.

Finance. International significance. The London economy was relatively fortunate in being able to offset manufacturing decline by participation in the growth of global financial markets. By 1990 one in six of London's workforce was in financial or business services—one-third of Britain's total employment in these sectors. The City claimed to have the largest concentration of financial employment in the world.

London's role as a world financial centre has long historical roots. At the end of the 19th century more than half the world's trade was financed in British pounds sterling. In the early 20th century the City played a more modest role as banker to the British Empire and the sterling area of trading nations. It regained a global presence thanks to the relaxation of exchange controls in 1958, the development of the Eurocurrency and Eurobond markets in the 1950s, and the deregulation of capital and securities markets—the "Big Bang"—in the 1980s.

In the 1990s London's most significant revenue-earning function was to provide a centre for the international banking market. A unique concentration of banks from every corner of the globe allowed an exceptional range of currencies to be traded. In the competitive climate of late 20th-century global finance, it captured a large share of activity in the newer, esoteric markets—swaps, cross-exchange equity trading, and currency options—and also maintained its historical position in the conventional fields of international bank lending, underwriting, bond trading, foreign exchange trading, and investment management. London had become Europe's main centre for large volume trading in securities. It dominated the world market in marine and aviation reinsurance.

Financial districts. The office towers of the financial services sector cluster tightly in the historic central business district of the City of London, or "Square Mile." Banking, insurance, maritime services, commodities, and stockbroking are each associated with particular districts within the City. The advantages of proximity and face-to-face dealing are offset by the shortage of space. The Square Mile is not Hong Kong or Manhattan or Chicago's Loop but a medieval city in which building opportunities are limited at every turn by the presence of ancient monuments and by fragmented and complex patterns of land ownership. During the 1980s, financial service activity began to spill into neighbouring areas with better availability of land. Obsolete railway stations provided many opportunities for office redevelopment—notably the renowned Broadgate project at Liverpool Street Station to the north of the Square Mile—as did the great premises in Fleet Street, to the west, that had been vacated by the newspaper industry in its shift from hot-lead production to computer typesetting. The most spectacular secondary centre was the business city of Canary Wharf, built by the Canadian Reichmann brothers in the derelict docks a mile and a half to the east of the Square Mile. The project bankrupted its developers but left London with an enduring memorial of the boom years of the financial services revolution.

Transportation. The Thames. London's oldest highway is the Thames. Until the opening of Westminster Bridge in 1750, London Bridge was the only crossing. Most journeys across the river and many trips within London were made by boat. Both banks were punctuated by stairs leading down to ferries. The Thames watermen, who had been regulated since the 14th century, formed their own guild or company by 1603. After the development of the railways, the river ceased to carry significant traffic in passengers, despite periodic attempts to revive its function as a mass transit artery with hydrofoils, catamarans, and hovercraft.

Roads. London's most striking physical feature is the absence of a grand road layout. Town planners have made repeated attempts to impose a greater degree of formal order on the capital. The most celebrated efforts in modern times have been Sir Patrick Abercrombie's Greater London Plan of 1944 and the Greater London Devel-

Production
of
consumer
goods

Centre of
the inter-
national
banking
market

Failure
of road
building

opment Plan of 1969, both of which attempted to drive modern highways through the fabric of the city. Fortunately, both plans were frustrated, leaving London with isolated stretches of high-speed road instead of a coherent network. A representative example is the quarter mile of dual carriageway that runs along the north of the Square Mile—a relic of Abercrombie's scheme for an inner ring road around the central business district. After years of economic blight, the scheme for a westward extension of the road was abandoned and the land sold off. A postmodern office block for a firm of city solicitors now closes the abandoned vista of the inner ring road known as Route 11.

The failure of road building has actually proved beneficial to London, which has Britain's highest rate of travel by public transport. The use of automobiles for travel to work in central London is small and declining. Each morning a million or so people enter central London, and well over three-quarters of them arrive by rail.

Railroads. The basis of the capital's rail infrastructure was laid in four heroic decades between 1836 and 1876. Competing railway companies brought 10 separate systems of track into London from every point of the compass, each with its own terminus station perched at the edge of the high-value metropolitan core of the City and West End. Linkage between the terminals was achieved in 1884 with the opening of the Metropolitan Railway, London's first "underground." Early development of underground railways in London was helped by the clay, which was easy to excavate, the spoil providing raw material to make bricks for lining the tunnel walls. Improved deep tunneling techniques after World War I allowed a rapid expansion of the underground network, while the Piccadilly, Bakerloo, Central, and Northern lines opened up hundreds of square miles of rural Middlesex and Essex for suburbanization. South of the Thames a similarly dense network of railway stations had been developed along the electric suburban lines of the Southern Railway out of Victoria, Waterloo, and London Bridge stations. Raised on handsome brick viaducts above the floodplain or sunk into cuttings through the rolling uplands, the converging railways form one of the most distinctive topographic features of south London. Together the underground of the north and the overground of the south equipped London with a network of lines and stations that is rivaled only by that of Tokyo for scale and density. Most of the network was already built and in service by 1939. In the second half of the 20th century it only remained to fill in the gaps in the network map with projects such as the Heathrow extension of the Piccadilly Line, the Jubilee Line, the Victoria Line, and the Docklands Light Railway.

Taxicabs. The black taxicab is a familiar feature of the London scene. The cabs and their drivers, the London cabbies, are products of a system of licensing that extends back to 1639. About 1900 there were more than 11,000 registered cabs plying for hire on the streets of London, and double that number a century later. Motorized cabs first made their appearance in 1904 and soon displaced the horse-drawn cab (the last survived until 1947). A regulation of 1906 required the new mechanical cabs to be designed in such a way that they could turn, like the two-wheeled hansom cab, in a circle of at most 25 feet. Still in force, it accounts for the surprising maneuverability of black cabs in congested London streets and their distinctive "sit up and beg" design. Cabbies themselves are subjected to even older regulations, which require them to pass a detailed test on topography, street names, and principal destinations throughout a six-mile radius from Charing Cross. Trainee taxi drivers acquiring "the knowledge," as it is called, are a familiar sight on the streets of the capital.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. *The City Corporation.* The Lord Mayor and City Corporation of London form one of the oldest local governments in the world, with a history of municipal autonomy extending in unbroken succession to the folk-moots of the early Middle Ages. The Square Mile remains an autonomous jurisdiction within its historic boundaries,

with its own police force and a complete range of municipal services attending to the needs of a resident population of less than 6,000 and a weekday working population of 300,000. Through the centuries the City Corporation has accumulated immense resources of capital and property, entrenching itself alongside its historic foes, the crown and the aristocracy, at the pinnacle of Britain's stratified society. The City has never concerned itself with the wider issues of local government in the metropolis, except insofar as they impinge upon its own position and privileges, which are tenaciously defended.

Envision/© Pictures of London



Modern and historic buildings cluster tightly in the City of London. A statue of the Duke of Wellington stands before the Royal Exchange (1844) at right, with the Bank of England (18th century) at left. The National Westminster Tower, at right rear, stands behind the two London Stock Exchange buildings at centre.

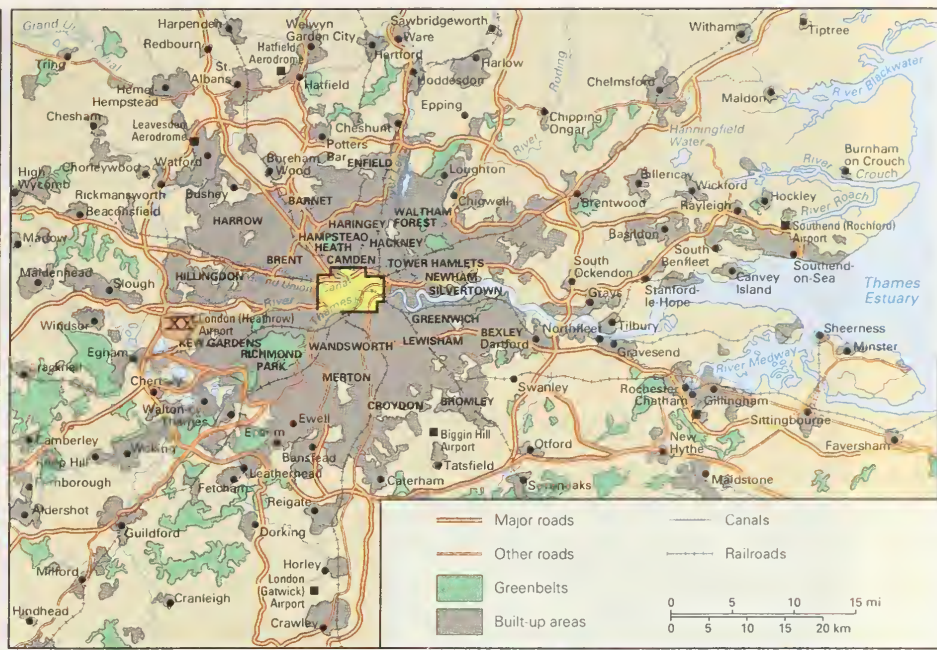
Around the unyielding nucleus of the Square Mile, arrangements for the wider metropolis have developed by stages. The City's indifference left a mid-19th-century population of almost three million under an anarchic miscellany of more or less undemocratic bodies based either on medieval ecclesiastical parishes or on ad hoc service agencies created under local legislation. In 1855 water and sewerage provision for the entire built-up area were brought under the control of the Metropolitan Board of Works. Following charges of corruption and lack of accountability, the organization was transformed in 1889 into the administrative nucleus of an elected local government for London as a whole, the London County Council (LCC). However, the City Corporation successfully lobbied to preserve its autonomy and secured the creation of a second tier of elected local governments, the metropolitan boroughs, to function as a political counterweight to the LCC.

Greater London. The same two-tier pattern, with its attendant tensions, was repeated in 1965 when the LCC was replaced by the Greater London Council. Its boundaries were extended to include all the suburbs developed after 1888—i.e., more or less the entire built-up area within

- | | |
|---|---------------------------------------|
| 1 Achilles Statue | 33 Old Admiralty |
| 2 Admiralty Arch | 34 Old St. Pancras Church Cemetery |
| 3 Bank of England Extension | 35 Parliament Square |
| 4 Bankside Power Station | 36 Polytechnic of the South Bank |
| 5 Buckingham Palace | 37 Queen Elizabeth Hall |
| 6 Central Criminal Court | 38 Queen's Gallery |
| 7 Covent Garden Market (site) | 39 Queen Victoria Memorial |
| 8 Downing Street | 40 Royal College of Art |
| 9 Euston Station | 41 Royal Courts of Justice |
| 10 Foreign & Commonwealth Office | 42 Royal Geographical Society |
| 11 Freemason's Hall | 43 Royal Opera House |
| 12 Geological Museum | 44 Somerset House |
| 13 Goods Depot | 45 Statue of Eros |
| 14 Government Offices | 46 St. Bartholomew's Hospital |
| 15 Hayward Gallery | 47 St. Clement Dane's Church |
| 16 Holborn Viaduct Station | 48 St. George's Catholic Cathedral |
| 17 Home Office | 49 St. James's Palace |
| 18 Horse Guards Parade | 50 St. Katherine Dock |
| 19 Imperial College of Science and Technology | 51 St. Margaret's Church |
| 20 Lancaster House | 52 St. Martin-in-the-Fields Church |
| 21 Lincoln's Inn | 53 St. Mary-le-Bow Church |
| 22 Lincoln's Inn Fields | 54 St. Mary-le-Strand Church |
| 23 London College of Printing | 55 St. Pancras Church |
| 24 London School of Hygiene and Tropical Medicine | 56 St. Pancras Hospital |
| 25 Ludgate Circus | 57 Temple Bar Memorial |
| 26 Marlborough House | 58 Trafalgar Square |
| 27 Ministry of Defence | 59 Treasury |
| 28 Museum of London | 60 Victoria and Albert Museum |
| 29 National Film Theatre | 61 Wellcome Museum of Medical Science |
| 30 National Portrait Gallery | 62 Wellington Arch |
| 31 National Temperance Hospital | 63 Wellington Museum |
| 32 Nelson's Column | 64 Westminster Abbey |
| | 65 Westminster Hall |



Central London and (inset) the Greater London metropolitan area.





Big Ben and the Houses of Parliament, as seen from the south bank of the Thames.

Envision/© V.K. Guy Photo Library

the Green Belt. At the same time, more than 100 existing local councils were amalgamated to form a modernized system of 33 boroughs, including the City of London and its young neighbour, the City of Westminster, chartered in 1900. The strategic authority and the local councils were at loggerheads from the outset, and their rivalries were intensified by shifts of political control at national, Greater London, and borough levels. Conservative Prime Minister Margaret Thatcher put these wranglings to an end in 1986 by the drastic expedient of abolishing the Greater London Council (then dominated by the opposition Labour Party) and stripping away its municipal assets. The majestic County Hall, across the bridge from the Palace of Westminster, was sold to a Japanese developer for conversion into a luxury hotel. London was left, uniquely among the world's great cities, with no form of civic leadership. Municipal affairs were distributed among the 33 borough councils, supplemented by a variety of ad hoc groupings and centrally appointed bodies for single services.

In a 1998 referendum, Londoners accepted a plan by the new Labour government to reestablish a citywide administration, which would operate in conjunction with existing Greater London boroughs. Established in 2000, the new Greater London Authority comprised a directly elected mayor and a 25-member assembly, and it assumed some of the local responsibilities that the central government had handled since 1986—notably over transport, planning, police, and other emergency services.

The London boroughs are large units, with populations usually of one or two hundred thousand, designed to achieve efficiencies of service provision rather than reflect local loyalty and sense of identity. When Londoners are asked where they live, they are more likely to give the name of their nearest railway station or shopping centre or of the pre-1965 administrative unit than of the modern borough.

Postal districts. At the local level, areas are often identified by postal district. The capital postal area is divided into 120 districts, each centred on a sorting office. An address in, say, SW1 carries a status that can be translated into property values. The dense urban core is divided into east central (EC) and west central (WC) areas with six numbered subdivisions (EC1–4 and WC1–2). The remainder of the metropolis is divided into six compass sectors, north (N), east (E), southeast (SE), southwest (SW), west

(W), and northwest (NW), with each sector comprising up to 28 numbered districts. The district nearest the city centre always has the number one, so that N1 (Islington), E1 (Stepney), SE1 (Southwark), SW1 (Westminster), W1 (Soho, Mayfair), and NW1 (Marylebone and Camden Town) form a ring around the central area.

The remaining London postal districts follow no geographic logic but are dotted randomly within their sector according to the alphabetical order of sorting-office names. Moreover, the outer boundaries of the London Postal Area fall short of the administrative and physical boundaries of the metropolis, giving several thousand suburbanites a postal address in the Home Counties though they live in a London borough. The one exception is the E4 postal district, which thrusts some miles out of London into the county of Hertfordshire.

Police. The Metropolitan Police was founded by Home Secretary Sir Robert Peel in 1829 and remains accountable to his successor, not to local councillors. By 1900 the Metropolitan Police District, which inherited responsibility for patrols against highwaymen, extended into the countryside in a 20-mile radius around London, an area of 700 square miles. Increased to 786 square miles, the area is still large enough to accommodate the entire metropolis and some of its rural fringe.

The Metropolitan Police District surrounds but excludes the City of London, which maintains its own police force. Metropolitan officers can be identified by their white shirts with silver buttons and City officers (who are recruited for height and wear lofty helmets) by their blue shirts with gold uniform trim. The British Transport Police and the Royal Parks Constabulary also maintain separate forces within the metropolitan area.

Hospitals. The history of London's great hospitals begins with medieval monastic charity. St. Bartholomew's, the oldest, was founded in 1123, and St. Thomas's at Lambeth dates to 1213. Other hospitals, all founded since 1700, include Guy's, St. George's, the London Hospital, the Middlesex, Charing Cross, the Royal Free, University College, and King's College. These hospitals combine the roles of medical school, research laboratory, general hospital, and specialized clinic at the highest level of excellence. Because of the concentration of medical resources in the heart of London, the administrative map of the health service divides the capital into four wedge-shaped regions, each hav-

The London boroughs

Two police forces

ing teaching hospitals at the inner tip and an outer catchment area extending far beyond the physical metropolis.

Education. School provision in London is a responsibility of the 33 boroughs. Nine out of 10 children attend borough schools. The remainder are at fee-paying private schools, of which the oldest and most august are Westminster School (originally monastic, refounded by Elizabeth I in 1560), St. Paul's School (1509), Harrow School (1572), Dulwich College (1618), and the City of London School (1834).

The panorama of higher education in London is characteristically complicated. Perhaps because of its civic fragmentation and the dominance of Oxford and Cambridge, the city lagged far behind other European capitals in the advancement of learning. The University of London, which was established as an examining body in 1836, did not become a teaching institution until 1900, centuries after its counterparts in Paris, Rome, and Madrid. Modern London has 12 universities in all, with more than 110,000 full-time and 50,000 part-time students. Despite the imposing monumentalism of its administrative buildings in Bloomsbury, the original London University is little more than a weak federation of 42 institutions ranging from small specialized schools to organizations such as Imperial College, University College, King's College, and the London School of Economics and Political Science, each of which operates in practice as a university in its own right.

Apart from a cluster of university buildings to the north of the British Museum in Bloomsbury, London's higher education facilities are spread widely through the metropolis. Halls of residence are even more scattered, and a high proportion of students live at home or in lodgings. The capital lacks an identifiable student quarter. Instead, that compound of offbeat bohemianism, nightlife, and political radicalism is sprinkled like yeast throughout Inner London.

CULTURAL LIFE

Centres of the arts. The competitive, localist streak that complicates public administration in London makes for exceptional cultural vitality. Artistic creativity flourishes in the diversity of rival centres of patronage. Royal patronage created Albert Hall, which every summer provides the setting for one of the world's greatest music festivals, the Sir Henry Wood Promenade Concerts known popularly as the "Proms." Municipal patronage, first of the London County Council and later of the Greater London Council, turned former industrial and warehousing land on the Waterloo riverbank into the South Bank arts complex, which combines the Royal Festival Hall, the Queen Elizabeth Hall, the Hayward Gallery, the Purcell Room, the National Film Theatre, the Museum of the Moving Image, and the Royal National Theatre. Not to be outdone, the City Corporation launched its own arts complex within the Square Mile at the Barbican, a high-density urban renewal scheme built on World War II bomb sites immediately north of the central business district. The Barbican has a concert hall, cinemas, an art gallery, a library, and a theatre that is the London home of the Royal Shakespeare Company. Each centre generates its own program of festivals and special events, as do borough councils and commercial promoters. No other city in Europe offers so many entry points to the young and talented musician, writer, artist, filmmaker, or performer. Though the figures are elusive, one estimate puts London's share of total national employment in cultural industries at 40 percent. Listings for the performing arts present a choice of more than 100 venues on a typical Friday or Saturday evening. Though the fragmentation of arts funding is often contrasted unfavourably with strong public sponsorship elsewhere, it is hard to resist the conclusion that London thrives on its distinctive combination of wide-open internationalism and local particularism.

Museums. The British Museum originated in 1753 in the government's purchase and amalgamation of three collections: the antiquities and natural history specimens assembled by the physician Sir Hans Sloane, the Cottonian Library and antiquities accumulated over 50 years by the Cotton family of Westminster, and the Harleian Collection

of Manuscripts built up by the 1st and 2nd Earls of Oxford. A public lottery raised the purchase price of the collections and a building in Bloomsbury to house them. This original nucleus was rapidly expanded by purchases and gifts as well as by the plunder of war and colonial conquest. In 1823-46 the Bloomsbury premises were totally rebuilt to the design of Robert Smirke, who graced the south front of the museum with a massive Ionic portico. The heart of Smirke's design, a large internal quadrangle, was roofed over in the 1850s with an immense copper dome to create the famous Reading Room in which Karl Marx wrote *Das Kapital*. In the late 1990s the inner courtyard and the Reading Room were enclosed by a 2-acre (0.8-hectare) square glass roof, transforming this area into what has been billed as the largest covered public square in Europe. Christened the Queen Elizabeth II Great Court, it was formally opened in December 2000.

The collections continued to outgrow the space available at Bloomsbury. During the 1880s the British Museum's botanical and zoological materials were reestablished in South Kensington as The Natural History Museum, housed in a richly carved and ornamented Victorian Romanesque building by Alfred Waterhouse. It formed part of a precinct of science and art developed at the direct initiative of Queen Victoria's husband, Prince Albert, on 87 acres of land purchased from the profits of the 1851 Great Exhibition. Its immediate neighbours are the Science Museum, the Geological Museum, and the immense collection of fine and applied arts from all over the world gathered under the roof of the Victoria and Albert Museum.

The full list of London museums has more than 250 entries catering to almost every industry, religion, ethnicity, profession, or foible. Opened in 2001, Firepower! The Royal Artillery Museum, at Woolwich, consolidates the holdings of the former Royal Artillery Regimental Museum and the Museum of Artillery at the Rotunda. Other outstanding collections include the National Maritime Museum at Greenwich (ships and the sea), the Museum of London at the Barbican (local history), the Museum of Mankind off Piccadilly (ethnography), the Bethnal Green Museum of Childhood (toys), and the Geffrye Museum in Shoreditch (domestic interiors).

Art galleries. *Exhibition spaces.* London is thought to possess about a third of the nation's art galleries and perhaps half the total hanging space in Britain. The greatest of the permanent collections is the National Gallery in Trafalgar Square. Behind it sits the National Portrait Gallery, which houses a vast collection of paintings, drawings, sculptures, etchings, photographs, and miniatures of famous faces past and present. The Tate's impressive holdings are displayed at two London locations: Tate Britain, at Millbank, which exhibits British art; and Tate Modern, at Bankside, where international modern painting and sculpture are housed. (The Tate also has galleries in St. Ives and Liverpool.) The Courtauld Institute Galleries at Somerset House on the Strand contain a major collection of French Impressionists and Postimpressionists. The Wallace Collection in Manchester Square combines paintings by great masters from several countries with furniture, ceramics, and goldsmiths' work in the ambience of an aristocratic town house. South of the river, Dulwich College has England's oldest public art gallery, built for important collections of 17th- and 18th-century masterpieces, including works by Rembrandt, Peter Paul Rubens, Thomas Gainsborough, and Nicolas Poussin. Its architect, Sir John Soane, incidentally also designed his own home in Lincoln's Inn Fields in 1812-13 to house an extraordinary personal collection of art (especially engravings and paintings by William Hogarth and Canaletto), works of art from the Middle Ages, and antiquities (classical and Egyptian). He bequeathed it to the nation, and today's visitors find it as he left it.

Both the National Gallery and the Tate mount special exhibitions. The other main venues for art shows are the Hayward Gallery on the South Bank, a sculptural concrete box of 1960s vintage, and the neoclassical Royal Academy of Arts in Burlington House on Piccadilly. The leading commercial galleries are concentrated in the West

Rival
centres of
patronage

The British
Museum

Com-
mercial
galleries

End of London around the epicentre of Bond Street. Specialist and avant-garde galleries are scattered throughout London, with a preponderance to the north and west.

Artists' quarters. Artists have long since been priced out of their traditional quarters in Chelsea and Hampstead. Today the prime areas for the bohemian life are the inner industrial suburbs of Hackney and Tower Hamlets to the east and Southwark to the south, where derelict work spaces and plentiful, loosely managed public rental accommodations have attracted many hundreds of artists. Scores of studios are thrown open in conjunction with the biennial exhibit of the Whitechapel Art Gallery in the East End.

Theatres. London offers every shade of dramatic experience, from authentic open-air Shakespeare performances in the replica Globe Theatre built by Sam Wanamaker to the bizarre offerings of the biennial London International Festival of Theatre (LIFT). On any given day the publicly sponsored National Theatre Company on the South Bank and the Royal Shakespeare Company in the Barbican have several shows in repertory. The West End has about 40 commercial theatres, playing to audiences composed in almost equal parts of Londoners, out-of-town theatregoers (many of whom come by coach), and overseas tourists. Many of London's centres—for example, Stratford, Ilford, Greenwich, Croydon, Battersea, Richmond, Kilburn, and Hammersmith—have fine theatres catering to local audiences. If one adds the fringe and studio performance venues, many in the upper rooms of pubs, the total offering of drama exceeds 100 shows a night before any count is made of amateur theatrical activity in churches, schools, and parish halls.

Music. The competitive ethos generated by London's administrative fragmentation is most evident in the realm of classical music. Five full-scale symphony orchestras vie for audiences and funding: the London Symphony Orchestra, the Royal Philharmonic Orchestra, the Philharmonia Orchestra, the London Philharmonic Orchestra, and the BBC Symphony Orchestra. On a slightly smaller scale, London also has the ensembles of the Academy of St. Martin-in-the-Fields, the London Sinfonietta, the City of London Sinfonia, the Sinfonietta 21, the Orchestra of the Age of Enlightenment, and the pit orchestras of the Royal Opera House in Bow Street, Covent Garden, and the English National Opera at the Coliseum Theatre north of Trafalgar Square. This immense pool of instrumental talent continually generates new performing groups and chamber ensembles. The Musicians' Union estimates that as many as 44 percent of Britain's working musicians are based in the capital.

Tributes. Every London district has a deposit of historical associations many centuries thick. Earlier generations of Londoners are present in street names, public statues and busts, and thousands of funerary monuments and inscriptions. Since 1867 London has paid tribute to distinguished citizens and visitors by attaching circular blue plaques to their places of residence. Among those honoured are Geoffrey Chaucer, Wolfgang Amadeus Mozart, Florence Nightingale, and James Joyce. The first plaque was put up on the wall of Lord Byron's birthplace in Holles Street, Westminster.

Sports. *Football.* Football (soccer) tops the lists of both participant and spectator sports in London. Amateur players turn out in the hundreds for games in every park and open space. Hackney Marsh, along the River Lea in the east of London, has a swathe of 100 pitches. The professional game, like almost every branch of London life, is organized locally, not citywide. As a result, nobody plays for London, yet the capital has 13 football clubs: Arsenal (based in Islington), Barnet, Brentford, Charlton Athletic, Chelsea, Crystal Palace, Fulham, Millwall, Orient, Queens Park Rangers (based in Shepherd's Bush), Tottenham Hotspur, West Ham United, and Wimbledon. Between them the clubs cover every part of London, and their colours can inspire strong local loyalty. The playing season lasts from August until May. Matches attract an average crowd of 15,000, rising to 30,000–40,000 for a big first-division game.

Cricket. In the summer months, county cricket and (in-

ternational) Test Matches are played at Lord's in St. John's Wood and at the Oval ground in Kennington, between Lambeth and Vauxhall on the south bank. The Surrey County Cricket Club has leased the Oval ground from the Duchy of Cornwall since 1845. London, England's most populous county, has no cricket team of its own but is partly represented by the historic counties corresponding to areas of the modern metropolis—Kent, Surrey, Essex, and the former Middlesex.

London has nearly 1,000 cricket clubs, and the amateur game is widely played on summer weekends, often on open greens and parks. The club at Woodford Green in north-eastern London claims to keep up the country's longest tradition of village cricket, beneath the statue of the former local member of Parliament Sir Winston Churchill.

Other spectator sports. June brings international tennis stars to the All England Lawn Tennis and Croquet Club at Wimbledon in southern London. An earlier highspot of the sporting calendar is the spring boat race between the Universities of Oxford and Cambridge, rowed up the turbulent waters of the tideway from Putney Bridge to Chiswick. Since the closure of the racecourse at Alexandra Park in September 1970, Londoners must travel out of town for the horse races—as they do by the thousands in June for the Derby on Epsom Downs, also in June for the Royal Week at Ascot near Windsor, and in July for the Goodwood races in Sussex. (M.J.H.)

Horse races

History

THE EARLY PERIOD

Foundation and early settlement. Although excavations west of London have revealed the remains of circular huts dating from before 2000 BC, the history of the city begins effectively with the Romans. Beginning their occupation of Britain under Emperor Claudius in AD 43, the Roman armies soon gained control of much of the southeast of Britain. At a point just north of the marshy valley of the Thames, where two low hills were sited, they established Londinium, with a bridge giving access from land to the south. The first definite mention of London refers to the year AD 60 and occurs in the work of the Roman historian Tacitus, who wrote of a celebrated centre of commerce, filled with traders. In the same year, Iceni tribesmen under Queen Boudicca (Boadicea) sacked the settlement. From traces of the fires they set, it can be determined that the city had already begun to spread across the Walbrook valley toward the hill where St. Paul's Cathedral was later built. After the sack, the city was reconstructed, including a great basilica—an aisled hall 500 feet long. On the same spot today stands Leadenhall Market, an 1881 creation of cast iron and glass. To protect the city, Cripplegate Fort was built by the end of the 1st century, with an amphitheatre nearby. The first half of the 2nd century was a prosperous time, but the fortunes of Londinium changed about AD 150, and areas of housing and workshops were demolished. A landward wall was built about AD 200 for defense. Remains of the wall can be seen at the edge of the Barbican (near the street called London Wall) and on Tower Hill. In medieval times the walls were rebuilt and extended, requiring new gateways in addition to the six Roman ones. During the 3rd century timber quays along the Thames and public buildings were rebuilt, and a riverside wall was constructed. An area of some 330 acres (about 135 hectares) was enclosed. Londinium in the 3rd and 4th centuries was less populous than in AD 125. When the legions were recalled to Rome early in the 5th century, the widespread abandonment of property. What happened to London over the next two centuries is a matter of conjecture.

The Roman foundation

Soccer



Londinium, c. AD 200.

a large and apparently densely built-up settlement (at least 150 acres) of craftsmen and traders just upstream of the depopulated Roman city and extending inland to what is now Trafalgar Square. The settlement was called Lundenwic; however, virtually nothing is known about this phase of London's history until the time of Alfred the Great (849–899) and the wars with the Danes, who invaded England in 865. A little farther west a church was founded on marshy Thorney Island in 785, later to be replaced by a great abbey (the Westminster) built at the behest of the pious Anglo-Saxon king Edward the Confessor.

Medieval London. The city's future importance as a centre of financial and military—and therefore political—power became clear at the time of the Norman Conquest (1066). One of the first acts of William I the Conqueror was to accord a charter promising the citizens of London that they should enjoy the same laws as under Edward the Confessor and that he would suffer no one to do them wrong. Just outside the city walls he established the Norman keep (the White Tower), which was the central stronghold of the fortress-castle known as the Tower of London. A roughly square (118 by 107 feet) structure, the White Tower is 90 feet high, with a tower at each corner of the walls. When in the late 12th century King Richard I returned from the Third Crusade with a new concept of fortification, he began surrounding the keep with concentric systems of curtain walls with towers at intervals, a project completed by Henry III (ruled 1216–72). Because virtually every reign since then has added its contribution, the Tower incorporates architecture from many periods. An official royal residence through the reign of James I in the early 17th century, it has also housed the Royal Mint, the Royal Menagerie, the public records, an observatory, an arsenal, and a prison. Some executions took place within the confines of the Tower, but most were carried out on Tower Hill just beyond. The Crown Jewels are now on display in the Tower, as is a superb collection of arms and armour.

The Norman kings selected Westminster as the site for their permanent residence and government. Edward the Confessor (ruled 1042–66) constructed an enormous church dedicated to St. Peter (and later referred to as Westminster Abbey) as well as a royal palace. The ancient

"city" of London, meanwhile, reestablished its role as a centre of trade. In 1085 London had between 10,000 and 15,000 inhabitants (less than 2 percent of England's population) and was the largest city in Europe north of the Alps. About 1087 a major fire destroyed many of the city's wooden houses and St. Paul's. In the rebuilding, houses of stone and tile began to appear, and some streets were partially cleansed by introducing open sewers and conduits, but wooden houses remained the norm. By 1200 the city and its suburbs involved a jurisdiction covering 680 acres (about 275 hectares)—which still defines the official limit of the City of London—and contained a population of 30,000 people. Between 1050 and 1300 construction of quays on the northern banks of the Thames led to the waterfront being extended southward by some 100 yards (90 metres). A colony of Danish merchants was outnumbered by Germans, who had their own trading enclave, the Hanseatic Steelyard, on the waterfront until they were expelled in 1598. Other important trading groups, who assimilated easily into London's population, were the Gascons, Flemish, and northern Italians. When members of the last group were firmly established as bankers, the Jews, who had arrived with the Normans, were banished in 1290; they were not to return until 1656.

In 1300 London had about 80,000 inhabitants that were provisioned by a food-supply network extending 40–60 miles into the surrounding countryside. The city also drew "sea coal" from Newcastle-upon-Tyne (300 miles distant by sea), and air pollution became a problem in London. The dynamism of this period came to a sudden end with the outbreak of the Black Death in 1348–49, with 10,000 Londoners being buried beyond the city walls at West Smithfield. Recovery of urban life was to prove a slow process.

By astute purchase from needy monarchs, the guilds—100 of them by 1400—were able to buy increasing freedom from royal intrusion in their affairs and further their self-government. The first mayor of London, Henry Fitz-ailwyn, probably took office in 1192. The first evidence of a Court of Common Council dates from 1332. Since disorder in the realm provoked unrest in the city, London usually supported strong, orderly government, especially in such crises as the deposition of Edward II (1327) and



The Royal Naval College on the Thames at Greenwich. Begun in the 1660s as a palace for Charles II, it was completed as a naval hospital from designs by Sir Christopher Wren; it became the Royal Naval College in 1873. At centre, between Wren's twin domes, is the Queen's House (now part of the National Maritime Museum).

© Denis Waugh

Richard II (1399), the Peasants' Revolt in 1381, and the rebellion headed by Jack Cade (1450).

Tudor London. By 1520 London was again enjoying prosperity, with 41 halls of craft guilds symbolizing that well-being. Toward the middle of the 16th century London underwent an important growth in trade, which was boosted by the establishment of monopolies such as those held by the Muscovy Company (1555), the Turkey (later Levant) Company (1581), and the East India Company (1600). It also grew in population, with the number of Londoners increasing from over 100,000 in 1550 to about 200,000 in 1600. The additional population at first found living space in the grounds of the religious institutions seized during the Reformation by Henry VIII (after 1536). To fill the void left by the cessation of the religious charities, the city organized poor relief in 1547, providing grain in times of scarcity and promoting the foundation or re-constitution of the five royal hospitals: St. Bartholomew's, Christ's, Bethlehem (the madhouse known as Bedlam), St. Thomas's, and Bridewell. Many of the private charities founded at this time are still in operation.

The population of the City and its surrounding settlements had reached 220,000 by the early years of the 17th century despite laws that attempted to contain the size of the capital. Indeed, the City Fathers (members of the Court of Common Council), tried to stop the subdivision of old houses into smaller, densely packed dwellings (a process known as "pestering"). New industries, including silk weaving and the production of glass and majolica pottery, were established, often outside the gates in order to avoid the restrictive regulations of the livery companies, which were successors of the craft guilds and were so named because of the distinctive clothing of their members. Slaughterhouses and numerous polluting industries were sited beyond the walls, especially to the east. The establishment of Henry VIII's naval dockyard at Deptford on the south bank was accompanied by a straggle of waterfront hovels on the north bank at Wapping.

When Henry VIII in 1529 began to convert Cardinal Wolsey's York Place into the royal palace of Whitehall and to build St. James's Palace across the fields, the City of Westminster began to take more definite shape around the court. Between Westminster and the City of London the great houses of nobles began to be built, with gardens down to the river and each with its own water gate. Along the Strand opposite these houses were distinguished lodgings for gentlemen who were in town during legal sittings. By the early 17th century the name London began to embrace both the City of London and the City of Westminster as well as the built-up land between them, but the two never merged into a single municipality.

The reign of Elizabeth I (1558–1603) arguably marked the apogee of the city's domination of England. The queen based her strength on its militia, its money, and its love. It provided one-quarter of the men for service abroad in 1585 and formed its armed "trainbands" (trained bands) to defend England against the threatened Spanish invasion.

17th-century London. The trainbands remained a force to be reckoned with, and Charles I, who had damaged the City's trading interests and flouted its privileges as cavalierly as he had Parliament's, was deterred from attacking London in 1642 by their presence at Turnham Green. Hostility toward the king made the fortified City the core of parliamentary support, and Parliament's success in the Civil Wars was due in good part to City allegiance.

In the early 1630s the 4th Earl of Bedford began developing Covent Garden, originally the convent garden of the Benedictines of Westminster, thereby initiating the process of building estates of town houses on land acquired from former religious houses.

In 1664–65 the plague, a frequent invader since the Black Death of 1348, killed about 70,000 Londoners. In 1666 the Great Fire of London burned from September 2 to September 5 and consumed five-sixths of the City. St. Paul's Cathedral, 87 parish churches, and at least 13,000 dwellings were destroyed, but there were only a few human fatalities. Because reconstruction had to be undertaken rapidly, adoption of a rational street plan was rejected, but the old streets were made wider and a bit straighter. Between 1667 and 1671 most of the houses were rebuilt (in brick since half-timbering was no longer allowed). Because many of the tiny parishes were combined and a few churches had escaped the fire, only about 50 churches were rebuilt, in addition to a new St. Paul's. Sir Christopher Wren, mathematician, astronomer, and physicist, though only informally trained as an architect, was given the formidable task of designing them and supervising their construction.

There is a famous inscription by Wren's son in St. Paul's Cathedral, addressing the visitor in the following words: *Lector, si monumentum requiris, circumspice* ("Reader, if you seek a memorial, look about you"). Much of the historic legacy of the City is in fact Wren's memorial. His churches are a series of virtuoso variations on basic architectural concepts. They range in style from the homely Dutch to the Gothic, but most of them embody his own conception of the classical style. The dome of St. Paul's is one of the most perfect in the world and, like the rest of the cathedral, is classical in theme with Baroque grace notes. The Monument for the Great Fire was adapted from a Wren design and erected near Pudding Lane, where the fire had started in the house of the king's baker. Wren constructed four other churches outside the City, built the Royal Hospital located in Chelsea, and designed parts of Kensington Palace, Greenwich Hospital, the Royal Observatory of Greenwich, and Hampton Court Palace.

Under Charles II royal abrogation of City rights was resumed, and, although James II restored forfeited City charters before his flight to France in 1688, it was in Guildhall under protection of the trainbands that the lords spiritual and temporal met to declare allegiance to William, the Dutch Prince of Orange (thenceforth known as William III of Great Britain).

To support the War of the Grand Alliance (1689–1713),

Relief for
the poor

The Great
Fire

Aspects of
the City's
power

City merchants in 1694 formed the Bank of England, and thenceforth the City's money market became a prime factor in the affairs of state. Another aspect of the City's power in the nation was the centring of the national press in Fleet Street (*The Times*, founded in 1785 off Blackfriars Lane, moved to new premises only in 1974). Finance, commerce, and port activities dominated the City and eastern London, while expansion of government and the attractions of fashionable society stimulated development of the West End.

As London continued to grow, the greater part of the metropolis lay outside the boundaries of the City. Whereas in 1550 75 percent of Londoners had lived under the Lord Mayor's jurisdiction, by 1700 (when there were 500,000 Londoners) only 25 percent did so, and in 1800 (1,110,000) the proportion was only 10 percent. Starting with Westminster Bridge (1750), half a dozen new bridges were built over the Thames, allowing new areas to be built up to the south. Important expansion occurred around the docks to the east as well as to the north and in the fashionable west. The rapidly expanding capital was governed by a patchwork of authorities, some of which were very ineffective. By 1700 London had overtaken Paris in population.

EVOLUTION OF THE MODERN CITY

18th-century London. By 1820, when George IV succeeded to the throne, many of the villages and hamlets that in the 17th and 18th centuries had been the destination of summer outings from the heart of the city had been covered by a tide of bricks and mortar. Some of the building was the well-planned work of great landowners; some, however, was the sorry work of the small or greedy. The Bedford, Portman, and Grosvenor estates, laid out on land that had passed from the monasteries into the hands of noble families, produced streets and squares that embellished the western part of town. On the other hand, to the east, parts of Stepney and Bethnal Green were constructed with ill-built cottage terraces. Agar Town and Somers Town, which lay near the modern King's Cross and St. Pancras railway stations, were very poorly built.

The changes brought during the years 1689–1820 followed no conscious plan. The government of the City was in full control and reasonably active within its jurisdiction. Beyond its boundaries, unchanged since the Middle Ages, government services and communications for the new areas came piecemeal. Important developers obtained local acts of Parliament enabling them to levy rates out of which to finance paving, lighting, cleansing, and the watch (a group of persons charged with protecting life and property). Because the popularity of the developers' streets depended in part on such services, they were usually adequately administered. Lesser developers left a legacy of slums and neglect for later generations to clear and repair.

Socially, commercially, and financially, London was the hub of the kingdom. It was also the centre of the world economy from the late 18th century to 1914, having taken over that role from Amsterdam. As a corollary to its great wealth, fed by the profits of the trade with the East and West Indies and with the Americas—indeed with most of the world—it reigned supreme in matters of the theatre, literature, and the arts. Eighteenth-century London was the city of David Garrick, Oliver Goldsmith, Samuel Johnson, and Sir Joshua Reynolds; of the great furniture makers and silversmiths; and of renowned foreign musicians, including George Frideric Handel, Joseph Haydn, and Mozart. London experienced important growth throughout the 19th century, with its total population exceeding 2,685,000 in 1851, the year of the Great Exhibition staged in Hyde Park to celebrate the commercial might of Britain and its empire. Fifty years later London's population reached 6,586,000, and the metropolis housed one-fifth of the population of England and Wales.

Organization, innovation, and reform. But the city's massive size brought increasing problems. Although new dispensaries and new or enlarged hospitals were reducing mortality, the former riverside town required new forms of government, communication, and sanitation if it was to continue to grow. These were slowly and painfully

introduced between 1820 and 1914, and the innovations came piecemeal. In 1829 a centralized Metropolitan Police force was provided, under the ultimate control of the home secretary, in place of the uncoordinated watchmen and parish constables. The lighting of streets by feeble oil lamps was revolutionized by the introduction of gas, and soon the Gas-Light and Coke Company (1812) was followed by similar companies scattered throughout London. Omnibuses (1829) began a revolution in passenger transport, and carriage by rail came less than 10 years later.

In 1842 an inquiry into public health exposed London's many deficiencies. Cholera in 1831–32 had caused the deaths of about 6,000 Londoners, and there were further outbreaks in 1848–49, 1854, and 1866. Legislation was passed in 1852 to assist provision of pure water. In 1854 the physician John Snow demonstrated the water transmission of cholera by analyzing water delivered by various private pumps in the Soho neighbourhood to a public pump well known as the Broad Street Pump in Golden Square. He arrested the further spread of the disease in London by removing the handle of the polluted pump. A statute of 1855 (the Metropolis Management Act) combined a number of smaller units of local government and replaced the medley of franchises with a straightforward system of votes by all ratepayers. Major works, such as main drainage and slum clearance, were put in the hands of the Metropolitan Board of Works.

The momentum of these changes, created by such reformers as Bishop C.J. Blomfield, Sir Robert Peel, Edwin (later Sir Edwin) Chadwick, and the Earl of Shaftesbury, continued throughout the century. New churches, new schools, better law and order, main drainage, pedestrian tunnels under the Thames, and care for social outcasts were some of the reformers' legacy. Their most visible bequests were Trafalgar Square, the Embankment, and roads, such as Shaftesbury Avenue and Charing Cross Road, driven through the worst of the slums. On another level, the School Board for London, established under the Education Act of 1870, set about the task of providing elementary education for all. Changes in local government continued, if not so drastically. The London County Council superseded the Metropolitan Board of Works in 1889, areas supervised by the vestries were reorganized into metropolitan boroughs by the London Government Act (1899), and various water companies were combined in 1902 into a publicly owned Metropolitan Water Board.

Public and private works continued to transform the appearance of London. The opening of the Metropolitan Line, a steam railway, in 1863 and the construction of Holborn Viaduct in 1869 were accompanied by the building of new Thames bridges and the rebuilding of Battersea, Westminster, and Blackfriars bridges. After years of discussion and agitation, the road bridges outside the City passed into public ownership, and the tollgates were removed. Main line railway termini were built on the edge of the built-up area (Paddington, Euston, St. Pancras, King's Cross), but eventually most of the railways from the south carried their lines across the Thames to the central business district on the north bank, with termini at Victoria, Charing Cross, Blackfriars, and Cannon Street. It was an era in which an abundance of initiative and capital was joined to abundant labour to make the widest use of new skills, cheap transport, and copious raw materials.

Technical progress continued gradually to alter the lives of Londoners and the appearance of the town. Cheap suburban train services enabled artisans or clerks to live farther and farther from their workplaces. Trams or streetcars (horse-drawn), after an unsuccessful beginning in 1861, became important in the 1870s and a major factor in metropolitan transport as their electrification developed in the early years of the 20th century. By then electricity was being used as the motive power for traffic below ground; the Prince of Wales (later Edward VII) opened the world's first electric underground railway, from King William Street in the City to Stockwell in the south, on Nov. 4, 1890. With the arrival, in 1897, of the gasoline-driven omnibus, transport in modern London was enhanced and the way opened for still faster development of suburbia. This was to be associated with the construction

Services for
London's
new areas

Railways
and bridges

of additional underground railway lines (especially north of the Thames) and the electrification of surface railway lines serving south London. The most celebrated suburban development involved what was to be known after 1915 as "Metroland" along the Metropolitan railway to the northwest of the capital.

Such changes were accompanied by rising land values in the central zone, by the construction of ever larger offices, factories, and warehouses in place of small houses, and by a continuous outlay of public and private funds on better housing and on street improvements. World War I, in which air raids inflicted some 2,300 casualties on London, brought only a temporary pause, and development resumed on a mounting scale after the war. The First British Empire Exhibition, held at Wembley in 1924-25, proclaimed the recovery of the nation and its colonies after the Great War. As a national and in some respects a world capital, London required institutions capable of meeting its needs. An era of amalgamation and expansion ensued, affecting almost all institutions. The docks continued to be enlarged, with London's last great enclosed dock (the King George V Dock) being opened in 1921. Street congestion increased, despite the rationalization of traffic authorities. By 1939 the population of the Greater London conurbation had risen to 8.6 million.

Reconstruction after World War II. London suffered widespread damage during World War II as a result of aerial bombardment, which devastated the docks and many industrial, residential, and commercial districts, including the historic heart of the City. About 30,000 Londoners died because of enemy action in the skies above the capital, and a further 50,000 were injured. The end of hostilities brought a return of evacuees, and reconstruction of the city began at once even though building materials were in desperately short supply. During the war the Greater London Plan (1944) had been prepared as a blueprint for reconstruction and also for relocating some Londoners and their jobs in new towns around the capital and in "assisted areas" in parts of the English provinces. Construction of new housing was discouraged and tightly controlled in a Green Belt around London, and the subsequent dispersed growth of the metropolis occurred in more distant sections of southeastern England. The New Towns Act (1946) gave rise to eight new settlements outside the metropolis. Passage of town and country planning acts, notably in 1947 and 1968, gave municipal authorities unprecedented powers of land purchase and control over development in London. The Festival of Britain (1951) proclaimed national recovery and produced the Royal Festival Hall on the south bank of the Thames as well as Lansbury Estate (a redevelopment area in Poplar). However, severe air pollution from coal-burning domestic hearths and industrial chimneys contributed to the great "smog" in 1952, which played a part in the death of 4,000 Londoners.

During the subsequent quarter century there was vast investment in slum clearance, construction of new houses and apartments, and improvement of services. Urban planning was more widely accepted, together with a broad policy to divert a share of employment and housing to localities beyond London's continuously built-up area. As a result, the number of residents in greater London contracted from about 8,193,000 in 1951 to about 6,600,000 in 1991; however, growth continued in outer parts of the southeast.

The port of London, which had been devastated during World War II, was restored in the 1950s. However, between 1968 and 1981 the city's docks were closed to traffic because of their small size, difficult labour relations and poor management, and powerful competition from major ports in continental Europe, especially Europoort in Rotterdam, Neth. During the 1980s the London Docklands Development Corporation encouraged major changes in Docklands, including construction of new housing and a large number of new offices (notably at Canary Wharf). London had experienced substantial deindustrialization by this time, with old industries that had been installed in Victorian times collapsing and many newer industries, dating from the interwar years and located along radiating

main roads laid out in those decades, sharing the same fate. London's economy had become increasingly geared to financial transactions and many other kinds of service activity. These sectors of the economy were strengthened by legislative changes in the mid-1980s affecting financial dealings. In consequence, the townscape of many parts of the City and the West End was transformed as vast new office complexes were constructed. Notable examples are Broadgate, on the site of the former Broad Street station, London Bridge City alongside the Thames, and the new Lloyd's building. In addition, London's airports at Heathrow and Gatwick were expanded, a major new airport opened at Stansted (30 miles north of the City), and a small airport for flights to western Europe began operating in Docklands. Completion of the M25 orbital motorway enabled vehicles to pass around the capital rather than move through it. However, road congestion remained a major problem, with even the M25 seriously overloaded with traffic. (B.E./H.D.C.)

New office complexes

Era of amalgamation and expansion

BIBLIOGRAPHY

Physical and human geography. BEN WEINREB and CHRISTOPHER HIBBERT, *London Encyclopaedia*, rev. and updated ed. (1993), is an indispensable, massive reference work with a detailed index; it covers locations, buildings, people, and historical events. KEITH HOGGART and DAVID GREEN, *London: A New Metropolitan Geography* (1991), collects academic essays on aspects of the society and economy of modern London. DONALD J. OLSEN, *Town Planning in London*, 2nd ed. (1982), a magisterial work, studies the interplay between aristocratic landowners and speculative builders that gave London its great 18th- and 19th-century estates. ALAN A. JACKSON, *Semi-Detached London: Suburban Development, Life, and Transport, 1900-39* (1973), provides a full and fascinating account of London's interwar suburbanization. GAVIN WEIGHTMAN and STEVE HUMPHRIES, *The Making of Modern London, 1815-1914* (1983), and *The Making of Modern London, 1914-1939* (1984), are good popular histories that are well illustrated. JOHN HILLABY, *John Hillaby's London* (1987), is an idiosyncratic essay on highways and byways, full of insights, by a celebrated "literary pedestrian." PETER HALL, *London 2001* (1989), provides a forward-looking analysis of change in the London region by England's best-known academic planner. (M.J.H.)

History. HUGH CLOUT (ed.), *The Times London History Atlas* (1991), contains more than 300 maps and illustrations in addition to substantial text and bibliography covering the development of the metropolis from its origins to the 1990s. ROY PORTER, *London: A Social History* (1994), provides an account of the capital from medieval times to the present. FELIX BARKER and PETER JACKSON, *London: 2,000 Years of a City and Its People* (1974, reissued 1984), provides a well-illustrated discussion of a broad time span. RALPH MERRIFIELD, *London: City of the Romans* (1983), presents recent archaeological discoveries. CHRISTOPHER N.L. BROOKE and GILLIAN KEIR, *London, 800-1216: The Shaping of a City* (1975), deals with political and economic changes in the early medieval period. MARY D. LOBEL and W.H. JOHNS (eds.), *The City of London from Prehistoric Times to c. 1520* (1989), contains scholarly essays and a number of extremely detailed maps. SYLVIA L. THRUPP, *The Merchant Class of Medieval London, 1300-1500* (1948, reprinted 1989), remains the classic scholarly study of the socioeconomic life of London's commercial elite. GERVASE ROSSER, *Medieval Westminster, 1200-1540* (1989), provides a detailed analysis of London's western settlement during its important medieval formative phase. JOHN SCHOFIELD, *The Building of London: From the Conquest to the Great Fire*, rev. ed. (1993), traces changes in the growth and internal character of the city over a wide period. JOHN STOW, *A Survey of London*, ed. by CHARLES LETHBRIDGE KINGSFORD, 2 vol. (1908, reprinted 1971), contains Stow's text of 1603 with additional notes tracing his sources of information. NORMAN BRETT-JAMES, *The Growth of Stuart London* (1935), is a well-referenced work based largely on original sources. T.F. REDDAWAY, *The Rebuilding of London After the Great Fire* (1940, reissued 1951), discusses the economic and social forces that shaped the rebuilding. M. DOROTHY GEORGE, *London Life in the XVIIIth Century* (1925, reissued 1984), chiefly records the life and work of poorer Londoners, using many quotations to provide contemporary points of view. JOHN SUMMERSON, *Georgian London*, new ed. (1988), a scholarly study, deals with the great estates and architecture of 18th-century London. DONALD J. OLSEN, *The Growth of Victorian London* (1976), well illustrated, discusses the key aspects of metropolitan growth during the 19th century. ANDREW SAINT (ed.), *Politics and the People of London: The London County Council, 1889-1965* (1989), contains thematic discussions of the important activities of the Council. (H.D.C.)

Los Angeles

A semitropical southern California metropolis of palm trees and swimming pools, television studios and aerospace factories, Los Angeles has become the second most populous city and metropolitan area (after New York) in the United States. The city sprawls across some 464 square miles (1,202 square kilometres) of a broad coastal plain agreeably situated between the San Gabriel Mountains on the east and the Pacific Ocean on the west. Its hallmark is an architecturally dramatic network of freeways. The automobile so dominates life in this uniquely mobile community that Reyner Banham, an English observer who took his cue from scholars who study Italian in order to read Dante, is said to have learned to drive a car so he could "read Los Angeles in the original."

The city is the seat of Los Angeles county, which contains more than 80 other incorporated cities, including Beverly Hills, Pasadena, and Long Beach, within its 4,081 square miles. The county also encompasses two channel islands, Santa Catalina and San Clemente; a mountain peak, Mount San Antonio, familiarly known as Old Baldy, 10,080 feet

(3,072 metres) high; more than 900 square miles of desert; and 75 miles (140 kilometres) of seacoast. The metropolitan area has paid for its spectacular growth by acquiring such urban attributes as smog-filled skies, polluted harbours, clogged freeways, crowded classrooms, explosive ethnic enclaves, and annual budgets teetering on the brink of bankruptcy. Since the city and the county are so intertwined physically and spiritually, any consideration of Los Angeles must move back and forth between the two entities.

Angelinos live in enclaves walled off from one another by ethnic, cultural, and economic differences. They go their own way, surfing, riding, skiing, yachting, hiking, playing golf and tennis. Nowhere in the world is the pursuit of happiness more unabashedly hedonistic, and, perhaps, no city in modern times has been so universally envied, imitated, ridiculed, and, because of what it may portend, feared. As early as 1927 it was recognized by Bruce Bliven as "a melting pot in which the civilization of the future may be seen bubbling darkly up in a foreshadowing brew."

This article is divided into the following sections:

Physical and human geography 299

The landscape 299

The city layout

Climate

Natural phenomena

The people 301

The economy 301

Agriculture

Services

Industry

Motion pictures

Transportation

The communications media

Administration and social conditions 302

Government

Education

Recreation and cultural life 302

Parks

Sports

Historic landmarks

The arts

History 303

The growth of the metropolis 303

The Spanish-Mexican town and city

The American city

The modern city 304

Bibliography 304

Physical and human geography

THE LANDSCAPE

The city layout. The city is grotesquely shaped, like a charred scrap of paper, with independent municipalities such as Beverly Hills and Culver City as well as unincorporated county land lying within its boundaries. Elevation averages about 275 feet, ranging from sea level to 5,082 feet at Mount Lukens (also called Sister Elsie Peak). The Santa Monica Mountains, covering an area of 92 square miles and reaching heights of 3,000 feet, bisect the city, separating Hollywood, Beverly Hills, and Pacific Palisades from the southern boundary of the San Fernando Valley (the Valley: home to "Valley girls"), a 220-square-mile area with such suburban communities as Burbank, Glendale, North Hollywood, Studio City, Sherman Oaks, Encino, Tarzana, Woodland Hills, and the mission city of San Fernando.

The Valley's principal east-west artery, Ventura Boulevard, is a 17-mile bazaar of specialty shops, ethnic restaurants, banks, medical buildings, shopping malls, automobile agencies, and realtors' offices. In the 1920s it was a dirt road. The post-World War II boom turned it into the main street of what would now be one of the country's largest cities if the Valley were an independent entity.

Once the sanctuary for middle-class white families fleeing the city's congestion and racial tensions, the Valley has broken up most of its rural estates to make room for condominiums and shopping centres. Walnut orchards and truck gardens have given way to housing for blacks, Hispanics, and Asians who have gone to work in new plants ranging from basic industry to high technology. Burbank, long the butt of television comedians ("Burbank has a low suicide rate, because living in Burbank makes suicide redundant"), now proclaims itself the country's entertain-

ment centre. It is home for recording companies, the National Broadcasting Company, and three major motion-picture studios (Walt Disney Productions, Warner Bros. Studios, and Columbia TriStar Television Division).

Hollywood, eight miles northwest of the central city, was laid out in 1887 by Horace Wilcox, a Prohibitionist, who intended his subdivision to be a sober, God-fearing community. In 1910, when its water supply ran low, Hollywood was gobbled up by Los Angeles. The following year Blondeau Tavern, at the intersection of Sunset Boulevard and Gower Street, was turned into Hollywood's first motion-picture studio—to be abandoned 60 years later when Columbia Pictures moved to Burbank. By then the stars had long since left Hollywood, many of them moving into secluded hillside mansions above Beverly Hills, the most famous of which was Pickfair, built by Douglas Fairbanks for Mary Pickford in 1919 (since demolished). This level of glamour has been eclipsed by such edifices as the 123-room mansion of television producer Aaron Spelling.

Climate. Coastal mountain ranges to the north and east act as buffers against extremes of summer heat and winter cold. Even in the hottest months, the humidity tends to be mercifully low and the nights cool. "Night and morning low clouds" is the most common summer forecast, with the sun breaking through in the afternoon. Pronounced climatic differences occur in different sections of the city. The San Fernando Valley is generally several degrees cooler in winter and warmer in summer than communities on the opposite side of the Santa Monica Mountains. The city's mean temperature is about 64° F (18° C). Daytime high temperatures can reach 90° F (32° C) in any month; on occasion nighttime lows approach freezing. The average annual rainfall is 14 inches (356 millimetres), with most of it falling in the winter months.



The city of Los Angeles and vicinity.

Natural phenomena. Smog. Juan Cabrillo, California's first European explorer, reported "many smokes" in 1542 hanging over "a large bay," which was probably the future harbour of Los Angeles. Four centuries later, in September 1943, Angelinos had their first massive doses of air pollution. Southern California's highly publicized sunshine, they learned, was creating an inversion, a warm layer of air that clamped a lid on the cooler air of the saucer-like basin below.

The sun continues to cook the noxious vapours rising from factories and oil refineries, and, most of all, from the tail pipes of the county's millions of cars, vans, and trucks. After many years and the expenditure of billions of dollars, Angelinos have had to settle for measures designed not to make the air really clean but to keep it from getting appreciably worse. No significant curtailment of smog can be achieved, they have been told by the head of their urban

air-pollution management agency, "without a major overhaul of the industrial and transportation structure, severe limits on growth, and a radical change in lifestyle."

Earthquakes. The great San Andreas Fault is 33 miles from downtown Los Angeles at its closest point, but more than 40 known lesser faults crisscross the metropolitan area. No earthquakes of major intensity—i.e., higher than 7 on the Richter scale of 10—have hit southern California since the 1850s, but, aside from the periodic minor jolts that Angelinos take for granted, destructive temblors have struck Santa Barbara (1925), Long Beach-Compton (1933), and the San Fernando Valley (1971 and 1994).

Winds. The city is occasionally buffeted by a Santa Ana, a hot, dry wind named for the canyon through which it often blows. Santa Anas occur when air rushes down from the high inland plateaus and is heated by compression.

Environmental blessings and blights

Angelinos are noticeably restless during a Santa Ana. The author Raymond Chandler wrote, "Meek little wives feel the edge of the carving knife and study their husbands' necks."

Fires. Dry winds whipping through narrow canyons on hot days heighten the ever-present danger of fire in the brush-covered mountains. Costly brushfires in 1961 destroyed 484 homes in Bel Air, an affluent west-side community; in 1993 brushfires throughout the Los Angeles region scorched 200,000 acres and destroyed more than 800 homes. Hillside homeowners run the risk not only of fire but also of mud slides when winter rains pelt canyon walls denuded by flames.

THE PEOPLE

When the city celebrated its bicentennial in 1981, whites had become a minority, as they had been when the city was founded. The 2000 census confirmed what everyday experience indicated, that Latinos had become the most numerous element of both the city and the county. Asians displaced African Americans as the third largest ethnic group. While some neighborhoods retained an ethnic solidarity, the clearest dividing lines were along economic strata.

Latinos. As throughout the state, Latinos increasingly demanded their share of jobs, opportunities, and political clout. Long confined to barrios where gangs are plentiful, streets drug-infested, schools overcrowded, and housing substandard, the community began to rise to middle-class status as it dispersed over the county. As whites and African Americans left the city centre, their places were filled by waves of Spanish-speaking immigrants. The language barrier, long a bar to social and economic progress, was becoming less of an impediment.

Asians and Pacific Islanders. In the late 1970s, the city's growing Asian and Pacific Island population helped to resettle refugees from Vietnam, Cambodia, and Laos. Little Tokyo, a few blocks from City Hall, has been the country's largest mainland concentration of Japanese Americans since the early 1900s. During World War II some 40,000 southern Californians of Japanese ancestry were placed in isolated detention camps. For two decades following their release, Little Tokyo languished, but in the 1980s it blossomed once again, redeveloped largely with capital from Japan. In both the city and suburbs, the distinctive Hangul signs announcing Korean businesses and restaurants could be seen.

African Americans. Although 26 of the city's 44 founders were of African ancestry and the granddaughter of a black tailor at one time owned the land now occupied by Beverly Hills, there were only 38,894 black Angelinos in 1931. Sixty years later, streams of blacks, mostly from the rural South, had raised the figure to about 500,000. By 2000, few solidly black communities remained in the central city, and the proportion of African Americans had dropped to about one-tenth of the overall population.

Watts, the 2.5-square-mile core of what was then the city's south-central black ghetto, exploded on a hot summer night in August 1965. The flames of Charcoal Alley gave whites a glimpse of the misery and anger of rural blacks trapped in an urban slum. After six days of burning and looting, 34 people were dead and 1,032 wounded and 3,952 had been arrested. Racial tensions erupted again into riots in 1992 after police were cleared of criminal charges in the beating of Rodney King; the episode had been caught on videotape and was endlessly replayed, fueling outrage. The community's interaction with police had long been a sore point, but the rage was turned against businesses and hapless civilians caught in the rampage. The riot resulted in 52 deaths and \$775 million in property damage.

THE ECONOMY

Once a land of vineyards, orange groves, and dairy farms, Los Angeles has become the nucleus of a vast industrialized urban area radiating from the city's bustling financial district. If this "Golden Circle" were a separate country, its gross national product would be among the highest in the world.

Agriculture. While its population rose almost 50 percent in the 1950s, Los Angeles county sacrificed 3,000 acres (1,200 hectares) of farmland a day to the bulldozer (freeways alone consumed about 40 acres per mile), and it ceased to be the nation's wealthiest agricultural county. Agriculture nonetheless continues to play a role in the county's economy; principal crops include nursery and greenhouse plants, vegetables, fruits, nuts, seeds, and hay.

Services. The service sector is the primary factor in the Los Angeles economy. Business and professional management services, health services and research, and finance are important, as are trade and tourism.

Industry. Before World War I, San Francisco was the state's manufacturing hub, but since the 1920s it has been far outstripped by Los Angeles. Major products include aerospace equipment (although cuts in federal defense spending reduced this industry in the 1990s), petroleum and refining, processed food, electronics, medical equipment, aluminum, furniture, automotive parts, and toys. Apparel design and production is a primary industry (and high-end shopping is a favourite pastime). High-technology industries—including the manufacture of computers, communications equipment, missiles, and space vehicles as well as computer software development—have become important to the modern economy.

Aviation. The first international air meet in the Western Hemisphere was held on the outskirts of Los Angeles in 1910. A dozen years later Donald W. Douglas was turning out an airplane a week in his Santa Monica plant. One-third of the nation's World War II military aircraft were built in the Los Angeles area. The city's preparation for the 1984 Summer Olympics included a major refurbishing of its international airport.

Oil. Los Angeles sits on one of its major resources, as Edward L. Doheny demonstrated in 1892 when he started selling local crude oil for industrial fuel. In 1921, after the automobile had created a ravenous appetite for petroleum products, the world's richest deposit in terms of barrels per acre turned up on Signal Hill. In its first half-century its 2,400 wells produced 859,000,000 barrels of oil.

James Randklev—Shostal Assoc.



The campus of the University of California at Los Angeles.

Motion pictures. Generations of moviegoers have looked up to Los Angeles as the world's film capital. It still plays a leading role in the motion-picture, television, radio, and recording industries, although some major studios have disposed of their valuable real-estate holdings. The back lot of the 20th Century-Fox Film Corporation was converted to Century City, a multimillion-dollar "instant city." The filmmaking activities at Universal Studios in the San Fernando Valley are a popular attraction for the area's millions of visitors each year.

Transportation. The vast majority of Angelinos working in the metropolitan area drive to their jobs in a car, van, or truck, while only a small percentage use public transportation. A light rail and subway system was begun in the late 1980s; its first segment opened in 1990. Ironically, much of the city's astonishing growth in the early 1900s was due to the superb interurban transit service provided by the big red electric cars of the railroad magnate and art collector Henry E. Huntington. The system became a freeway casualty beginning in the early 1940s.

The Port of Los Angeles is one of the busiest in the nation and the world. Located in San Pedro Bay, about 20 miles south of downtown, its international trade includes automobiles, furniture, apparel, paper products, chemicals, and grains. It is also an important cruise-ship port. Main trading partners are Japan and China. Los Angeles International Airport began as an airfield in 1928 and started commercial airline service in 1946.

The communications media. The city's first newspaper, the *Star*, began weekly publication in 1851. Thirty years later the *Los Angeles Times* published its first issue. Acquired the following year by Gen. Harrison Gray Otis, it became the bible for the city's boosters, conservative Republicans, and antilabour forces. When it passed into the hands of Otis' grandson, Otis Chandler, in the 1960s, it took a more liberal and worldly stance. After the demise in 1989 of its afternoon rival, the *Herald-Examiner*, the San Fernando Valley's *Daily News* became its chief competitor. More than two dozen of the area's radio stations broadcast in languages other than English. The Spanish-language television station KMEX holds its own against network stations.

ADMINISTRATION AND SOCIAL CONDITIONS

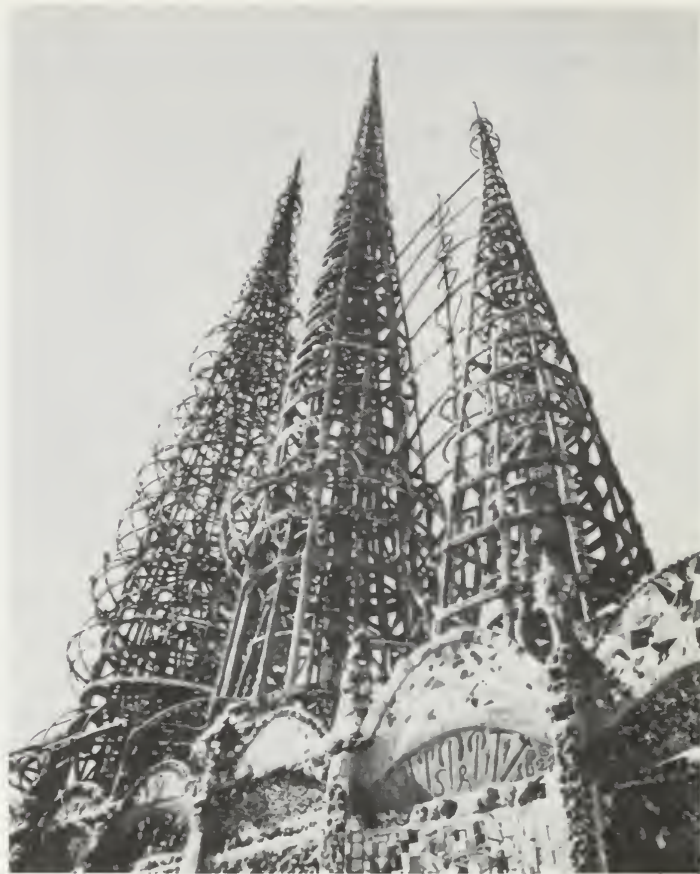
Government. The city's mayor and its 15 council members are elected to four-year terms, along with the city attorney, controller, and the seven-member board of education. The county is run by a five-member board of supervisors serving four-year terms. They preside over a jurisdictional jungle of overlapping city-county agencies.

Education. There are about 20 junior colleges in the county and five state colleges. The area's two oldest institutions of higher learning are the University of Southern California (USC; 1880) and Occidental College (1887). USC is noted especially for its schools of law, medicine, dentistry, engineering, and performing arts. The state-supported University of California at Los Angeles (UCLA; 1919) has a wide range of undergraduate and graduate offerings, with a particularly heavy commitment to the life and geophysical sciences and the arts. California Institute of Technology, founded in 1891 as Throop Polytechnic Institute, moved to its present location in Pasadena in 1910. Researchers there conduct probes of outer space in conjunction with the Jet Propulsion Laboratory and the Mount Wilson Observatory.

RECREATION AND CULTURAL LIFE

Parks. Hancock Park, the 23-acre site of the Rancho La Brea pits, was given to the county in 1916 by G. Allan Hancock, an oil magnate. Excavations between 1906 and 1913 of the black, bubbling pools had revealed a unique repository of fossilized skulls and bones of long-extinct mammals trapped in the seepage of *brea*, or "tar" (actually asphalt). In addition to the tar pits, the park features life-size figures of such creatures as the imperial mammoth, the American mastodon, the sabre-toothed tiger, the ground sloth, and the short-faced bear.

Griffith Park is spread across some six square miles of rugged mountains, an area larger than Beverly Hills. Mule



The Watts Towers by Simon Rodia, completed 1954, in Los Angeles.

Tom McHugh—Photo Researchers

deer wander down from the hillsides to peer at the exotic creatures in the city's 80-acre zoo. Within easy reach of the city by freeway are such commercial ventures as Disneyland, Knott's Berry Farm, and Magic Mountain.

Sports. Intercollegiate athletics is highlighted by the intense USC-UCLA rivalry. The annual New Year's Day Rose Bowl football game is played in Pasadena. Angelinos also support professional major-league teams in baseball, football, basketball, ice hockey, and association football (soccer). There is horse racing at Santa Anita Park and Hollywood Park.

Historic landmarks. Angelinos have held onto little of their Spanish-Mexican past, but largely through the efforts of Charles F. Lummis, a colourful turn-of-the-century editor and writer, the missions of San Gabriel Arcángel (1771) and San Fernando Rey de España (1797) have been preserved. Olvera Street, a narrow, block-long string of Mexican shops, cafés, and the Ávila family's adobe townhouse (c. 1818), has been a popular tourist attraction since its opening in 1930. It is now part of El Pueblo de Los Angeles Historic District, an area of some 40 acres, which includes the plaza where the city got its start and the first church its residents built (dedicated in 1822, rebuilt in 1861-62).

The three Towers of Simon Rodia (99, 97, and 55 feet high), better known as the Watts Towers, were built of broken tiles, dishes, bottles, and seashells over a 33-year period by Rodia, an unschooled Italian immigrant who later explained, "I had in mind to do something big, and I did." When his work was completed in 1954, he gave the property to a neighbour and left, never to return.

The city's early 20th-century Mission-style architecture gave way to the wood-shedded California bungalow modeled on the work of Charles and Henry Greene, whose 8,000-square-foot Gamble House, built in Pasadena in 1908, is now used as a study centre for student architects and designers. Frank Lloyd Wright received several commissions in southern California in the early 1920s following completion of Hollyhock House in what is now Barnsdall Art Park. R.M. Schindler helped supervise the

Demise
of the
urban
transit
system

Institu-
tions of
higher
learning

project before setting out to build his own masterpieces. Irving Gill is generally regarded as the city's most neglected architect. Richard Neutra's steel-frame Lovell House, built in the late 1920s, stands as a monument to the International Style; striking structures by Richard Meier and Frank Gehry set the tone for the last quarter of the 20th century.

One of the delights of downtown Los Angeles is the Bradbury Building (1893). Sunlight warms and illuminates the five-story inner court with its delicate French ironwork, Belgian marble, and Mexican tile. Union Station (1939), among the last of the country's railroad cathedrals, is a charming interpretation of Spanish Mission architecture. The city's Central Library (1926), incorporating Egyptian motifs, was the last building designed by Bertram Goodhue, whose low, buff-coloured, stuccoed building recalls the work of Irving Gill, although the basic form was Beaux-Arts. After two fires in 1986, a new wing was built, doubling the library space, and a striking eight-story atrium was added.

The arts. Los Angeles has from time to time sheltered outstanding composers and writers such as Igor Stravinsky, William Faulkner, F. Scott Fitzgerald, and Aldous Huxley. The city has been the subject of innumerable novels, the most durable of which appear to be Fitzgerald's *The Last Tycoon* (1941), Huxley's *After Many a Summer Dies the Swan* (1940), Evelyn Waugh's *The Loved One* (1948), and Nathanael West's *Day of the Locust* (1939). Crime novels set in Los Angeles abound. In addition to those by Raymond Chandler and Ross Macdonald, the Easy Rawlins series by Walter Mosley, set in Watts, is notable for its sense of time and milieu and its African American protagonist. The city is also well known to the world's book collectors for the works of such fine-printers as Saul Marks, Grant Dahlstrom, and Ward Ritchie.

Los Angeles has developed a lively marketplace for works of art. Many of its leading galleries are scattered among the fine restaurants, antique shops, and rare book stores on La Cienega and Melrose boulevards in West Hollywood.

Performing
arts

The performing arts are housed at the Music Center, which includes the elegant Dorothy Chandler Pavilion, home of the Los Angeles Philharmonic Orchestra; the Ahmanson Theater, used for plays, musical comedies, and light operas; and the cake-shaped Mark Taper Forum, designed for experimental stage productions. The Los Angeles Opera opened in 1985, filling a long-felt need. There is no resident ballet company, but visiting companies regularly perform at the Music Center, and Los Angeles-based companies present performances of modern, tap, jazz, ethnic, and experimental dance.

The Hollywood Bowl, a natural amphitheatre in the Hollywood Hills, offered its first production, *Julius Caesar*, on May 19, 1916, and initiated its programs known as Symphonies Under the Stars in 1922. The out-of-doors Greek Theatre in Griffith Park provides a wide range of summer entertainment, mostly musical. The Universal Amphitheatre in Universal City has offered year-round entertainment since it reopened in 1982 with a dome covering its seats. The Shubert Theatre in Century City is a few blocks from the massive 1880s New York set built for the motion-picture version of *Hello, Dolly!* (1969). Small, innovative theatres continue to spring up in the motion-picture capital.

Museums
and
libraries

The delightful California Science Center—an interactive facility, opened in 1998 in Exposition Park, not far from the Natural History Museum—is one of the largest natural history museums in the United States. Its unique fossil collection has been transferred to the Page Museum at La Brea Discoveries in Hancock Park. It shares the historic acreage with the Los Angeles County Museum of Art, which seems to float in the black pools of tar. Arata Isozaki, a Tokyo-based architect, designed the red sandstone, aluminum, and glass Museum of Contemporary Art (MOCA), housing art from the 1940s to the present. Large-scale installations are displayed in MOCA's Geffen Contemporary venue. The Southwest Museum in Highland Park, devoted to Native American art and artifacts, has one of the country's richest collections of textiles, pots, baskets, photographs, and books in this field. A touch-and-play Children's Museum has enriched and enlivened the Civic Center, and a Junior Arts Center is part of the Barnsdall Art Park complex built around the Municipal Art Gallery.

Pasadena's Norton Simon Museum houses a notable art collection spanning 2,000 years. In nearby San Marino the Henry E. Huntington Library, Art Collection, and Botanical Garden is a hospitable, parklike retreat, where visitors are invited to enjoy its roses, camellias, desert plants, and Japanese Garden as well as its wealth of rare books, manuscripts, and British art works. The travertine-clad Getty Center, designed by Richard Meier, opened in 1998 on a 110-acre hilltop campus. The Center, funded by a \$700,000,000 bequest from oil magnate J. Paul Getty, houses a world-class collection of fine and decorative arts and is a major tourist draw. In addition to local libraries in such cities as Pasadena, Santa Monica, Burbank, and Beverly Hills, Angelinos draw on the resources of the city's system and the county's regional and community libraries. UCLA has one of the country's finest academic libraries. The off-campus William Andrews Clark Memorial Library is noted for its John Dryden collection.

History

THE GROWTH OF THE METROPOLIS

The Spanish-Mexican town and city. On August 2, 1769, a Spanish expedition headed by Gaspar de Portolá, searching for mission sites, camped near a river they named the Porciúncula, in honour of Our Lady the Queen of the Angels of Porciúncula (Nuestra Señora la Reyna de los Angeles de Porciúncula). The area's first Europeans exchanged gifts with the peaceful, Uto-Aztecanspeaking peoples of the nearby village of Yang-na and left on the following morning. Despite three earthquakes during his overnight stay, Father Juan Crespi noted in his diary that "this delightful place among the trees on the river" had "all the requisites for a large settlement."

Early
settlers

Two years later the Mission San Gabriel Arcángel was established about nine miles northeast of the campsite. A decade went by before Gov. Felipe de Neve succeeded in colonizing the fertile river basin with 44 recruits from Mexico, most of them of Native American and African descent. The illiterate settlers assembled on the west bank of what is now the Los Angeles River on September 4, 1781, to claim the land they had been promised. Little is known about the events of that momentous day, but mythmakers have cloaked the city's founding in a ceremonial splendour worthy of its destiny and its high-sounding name, which has long been confused with the name given the river. It is now generally agreed that the city's correct name is El Pueblo de la Reyna de Los Angeles ("The Town of the Queen of the Angels").

El Pueblo, as it was commonly called, remained so isolated from the United States during the settlement's formative years that Joseph Chapman, the first Yankee to become an Angelino, was thought of as an Englishman ("El Inglés"). An engaging pirate from Boston, Chapman landed in 1818 with a black fellow privateer, Thomas Fisher. The first outsider to arrive by way of the arduous overland route was a fur trapper, Jedediah Smith, who turned up in 1826, four years after an independent Mexico had hoisted its flag above El Pueblo.

Los Angeles, with a population of nearly 1,250, had become a *ciudad* (city) in 1835 when Richard Henry Dana looked in on it. "In the hands of an enterprising people," he mused in *Two Years Before the Mast* (1840), "what a country this might be." By the time war broke out between the United States and Mexico in 1846, the California capital was so overrun with enterprising people that Gov. Pio Pico felt helpless against "the hordes of Yankee immigrants," who were "cultivating farms, establishing vineyards, erecting mills, sawing up lumber, building workshops, and doing a thousand other things which seem natural to them, but which Californians neglect or despise." When American forces under Capt. John C. Frémont and Commo. Robert F. Stockton entered the city on August 13, 1846, not a shot was fired. A revolt was put down the following January, and on July 4, 1847, Los Angeles celebrated its first Independence Day.

The American city. Los Angeles was incorporated on April 4, 1850, and designated the seat of Los Angeles county. The lawless, adobe cow town—"gambling, drinking and

whoring are the only occupations," grumbled a pioneer physician in 1849—prospered in the wake of the Gold Rush when hungry miners in San Francisco and Sacramento gorged on beef from southern California. A catastrophic drought (1862–65) following several years of declining cattle prices brought an end to the era of the ranchos. Vast Spanish and Mexican land grants, mortgaged and bankrupt, their owners ignorant of *Yanqui* laws and *Yanqui* interest rates, were broken up, fenced, and planted by a new breed of Angelino. By 1860 the city had become so Americanized that it had banned bullfighting and formed a baseball club.

With the arrival of the railroads (the Southern Pacific in 1876, the Santa Fe in 1885), Los Angeles began to ship its oranges back east and, by means of a massive advertising campaign, to lure immigrants westward to the New Eden. Aided by a railroad rate war, the boom of the 1880s more than quadrupled the city's population, from 11,183 in 1880 to 50,395 in 1890. Father Crespi's campsite was overrun by shrewd, aggressive Yankee boosters determined to build a city in their own image and requiring only two things denied them by a bountiful providence: a harbour and an adequate water supply.

The search for a port and water

Unlike San Francisco and San Diego, Los Angeles has no natural harbour. A narrow, artfully gerrymandered "shoestring strip" connects the inland city with its port, 23 miles south of City Hall. Los Angeles acquired its two harbour communities, San Pedro and Wilmington, by consolidation in 1909 following an epic struggle between the powerful Southern Pacific Railroad interests and the Free Harbor League. Work began on the harbour in 1899. The opening of the first municipal wharf in 1914 coincided with the completion of the Panama Canal, which put the ports of the Atlantic seaboard some 8,000 miles closer to Los Angeles. Its harbour became the busiest on the West Coast.

In 1904, casting about for new sources of water to sustain the city's relentless growth, William Mulholland, water-bureau superintendent, explored the Owens Valley some 250 miles northeast of Los Angeles and returned with a bold plan for an aqueduct to carry melted snow from the southern slopes of the Sierra Nevada to Los Angeles faucets. The plan outraged Owens Valley ranchers, enriched two syndicates of Los Angeles speculators, gave rise to rumours of wrongdoing and, highly fictionalized, ended up on the motion-picture screen as *Chinatown* (1974).

"There it is; take it," Mulholland told the thousands of Angelinos who assembled on November 5, 1913, to watch the Owens River water come cascading into a San Fernando Valley spillway. The 233-mile-long aqueduct, with 142 separate tunnels totaling 52 miles in length, has been supplemented by a 105-mile extension into the Mono Basin. The system supplies 80 percent of the city's water needs. The remainder comes from local wells, the California Aqueduct, and the Colorado River.

A severe drought in 1976–77 gave Angelinos a foretaste of the time (1985) when their Colorado River water would, by court order, be shared with Arizona. Meanwhile, with rain and snow dumping about three-fourths of California's usable surface water in the north, while more than half of its people live in the south, the "water wars" between powerful corporate farmers and resourceful urban environmentalists continue to be fought out in the legislature, the courts, and the voting booth.

THE MODERN CITY

In the first decade of the 20th century, while San Francisco tidied up the rubble of its 1906 earthquake, Los Angeles tripled its population, from about 100,000 to nearly 320,000. One local entrepreneur opened the first motion-picture theatre in the United States, in 1902, and another built the city's first garage to accommodate its growing number of automobiles, but the parasol-shaded girls from the red-plush brothels run by Pearl Morton and Cora Phillips still drove about in open carriages, much to the distress of the retired druggists, dentists, and wheat farmers from the Midwest who kept streaming into Los Angeles.

"Virtue has become virulent," reported a writer in *Smart*

Set (March 1913), who described an "overgrown village" swarming with "spiritualists, mediums, astrologists, phrenologists, palmists and all other breeds of esoteric wind-jammers." Successive waves of migratory writers, taking much the same tack, have stereotyped the city as an open-air institution for the eccentric.

The bizaare in the city

"Los Angeles represents the ultimate segregation of the unfit," Bertrand Russell declared, and Frank Lloyd Wright agreed. "It is as if you tipped the United States up and all the commonplace people slid down there into Southern California," he remarked in 1940, but within the next few years Los Angeles had become so overrun with such gifted European refugees as Heinrich and Thomas Mann, Arnold Schoenberg, Bertolt Brecht, Bruno Walter, Franz Werfel, Lion Feuchtwanger, and Alfred Neumann that the city had come to be dubbed "The Fourth Weimar Republic."

In the ensuing decades, Los Angeles prospered for its both natural and manmade attractions. When Disneyland opened in nearby Anaheim in 1955, it was an immediate hit and brought a new influx of tourism to the area. The 1984 Summer Olympic Games were a financial success, despite their being boycotted by communist nations in retaliation for the U.S. boycott of the 1980 Summer Games in Moscow. Los Angeles had also been host to the 1932 Summer Games.

But all is not perpetual sunshine: the city struggles with both natural and manmade disasters. An earthquake in January 1994, centred in the San Fernando Valley, resulted in some 60 deaths and the terrifying collapse of sections of major freeways. Forest fires and mud slides present a continual threat to multimillion-dollar houses. In 1999 and 2000 a police corruption scandal engulfed the city, and ethnic tensions are chronically close to a flash point.

BIBLIOGRAPHY. Introductions to the history of the city are JOHN CAUGHEY and LAREE CAUGHEY (eds.), *Los Angeles: Biography of a City* (1976); and JOHN D. WEAVER, *Los Angeles: The Enormous Village 1781–1981* (1980). LEONARD PITT and DALE PITT, *Los Angeles A to Z: An Encyclopedia of the City and County* (1997), is a useful and well-researched reference guide. Three standard works, although dated, are still indispensable: J.M. GUINN, *A History of California and an Extended History of Los Angeles and Environs*, 3 vol. (1915); CAREY MCWILLIAMS, *Southern California Country: An Island on the Land* (1946, reissued as *Southern California: An Island on the Land*, 1973); and WRITERS' PROGRAM, *Los Angeles: A Guide to the City and Its Environs*, completely rev. 2nd ed. (1951). The city's early years are discussed in ROBERT GLASS CLELAND, *The Cattle on a Thousand Hills: Southern California, 1850–1880*, 2nd ed. (1951, reissued 1990); and ROBERT M. FOGELSON and ROBERT FISHMAN, *The Fragmented Metropolis: Los Angeles 1850–1930* (1967; reissued 1993). Two volumes that focus on southern California and Los Angeles from 1850 through the 1920s are KEVIN STARR, *Inventing the Dream: California Through the Progressive Era* (1985), and *Material Dreams: Southern California Through the 1920s* (1990). Urban planning and the urban-suburban paradox are explored in GREG HISE, *Magnetic Los Angeles: Planning the Twentieth-Century Metropolis* (1997); WILLIAM FULTON, *The Reluctant Metropolis: The Politics of Urban Growth in Los Angeles* (1997); and ALLEN J. SCOTT and EDWARD W. SOJA (eds.), *The City: Los Angeles and Urban Theory at the End of the Twentieth Century* (1996).

Three early-day Angelinos have left invaluable recollections: HORACE BELL, *Reminiscences of a Ranger* (1881, reissued 1965–67), although it should be read with caution; HARRIS NEWMARK, *Sixty Years in Southern California, 1853–1913*, 4th ed. rev. and augmented (1970, reissued 1984); and BOYLE WORKMAN, *The City That Grew* (1936). One of the most useful records of pioneer days is JOHN ALBERT WILSON, *History of Los Angeles County, California* (1880, reprinted as *Reproduction of Thompson and West's History of Los Angeles County, California*, 1959). JOHN WALTON, *Western Times and Water Wars* (1992), is an excellent study of this key subject. The ethnic diversity of Los Angeles is explored in LYNELL GEORGE, *No Crystal Stair: African-Americans in the City of Angels* (1992); and ANTONIO RÍOS-BUSTAMANTE and PEDRO CASTILLO, *An Illustrated History of Mexican Los Angeles, 1781–1985* (1986). DAVID GEBHARD and ROBERT WINTER, *Architecture in Los Angeles* (1985), provides both description and history of important architecture. The first child of the freeway generation to make a serious appraisal of his off-ramp heritage is DAVID BRODSKY, *L.A. Freeway: An Appreciative Essay* (1981).

(J.D.W./Ed.)

Luther

The founder of the 16th-century Reformation and of Protestantism, Martin Luther is one of the pivotal figures of Western civilization, as well as of Christianity. By his actions and writings he precipitated a movement that was to yield not only one of the three major theological units of Christianity (along with Roman Catholicism and Eastern Orthodoxy) but was to be a seedbed for social, economic, and political thought. For further treatment of the historical context and consequences of Luther's work, see PROTESTANTISM.

By courtesy of the Nationalmuseum, Stockholm



Luther, oil painting by Lucas Cranach, 1526. In the Nationalmuseum, Stockholm.

LUTHER AS EDUCATOR AND MONK

Early life and education. Martin Luther was born on November 10, 1483, at Eisleben in Thuringian Saxony (Germany). His parents, Hans and Margarethe Luther, who had moved there from Möhra, soon moved on again to Mansfeld where Hans Luther worked in the copper mines, prospering enough to be able to rent several furnaces and to obtain a position among the councillors of the little town in 1491. Luther's few recollections of childhood that have survived reflect a sombre piety and strict discipline common in that age. His schooling seems to have been unremarkable: the Latin school at Mansfeld, a year at a school in Magdeburg (run by Brethren of the Common Life, a medieval lay group dedicated to Bible study and education) and at Eisenach in his 15th year, where he made valued older friends. In the spring of 1501 he matriculated in arts at the University of Erfurt, one of the oldest and best attended universities in Germany. There he talked long and seriously enough to be nicknamed "the Philosopher," and played the lute. He took the usual arts course and graduated with the B.A. degree in 1502. He took his M.A. in 1505, placing second among 17 candidates. In an age when few students got as far as the master of arts degree, he had fulfilled his parents' hopes. Like many other parents of his time, Hans Luther intended his son to become a lawyer, and he paid cheerfully enough for the expensive textbooks when Martin began legal studies. He was chagrined to learn that his son, without consulting his parents, had decided to enter religion and had sought admission to the house of Augustinian Hermits in Erfurt.

Brother Martin Luther. Evidence on the reason for his

decision to enter the religious life is scanty. In his later, not always reliable, *Tischreden* ("Table Talk"), it is related that on July 2, 1505, he was returning from a visit to his parents when he was overtaken by a thunderstorm near the village of Stotternheim and cried out in terror, "Help, St. Anne, and I'll become a monk." In his *De votis monasticis* ("Concerning Monastic Vows," 1521) Luther says "not freely or desirously did I become a monk, but walled around with the terror and agony of sudden death, I vowed a constrained and necessary vow." He sold most of his books, keeping back his Virgil and Plautus, and on July 17, 1505, entered the monastery at Erfurt.

Augustinian Order at Erfurt. In joining the eremitical order of St. Augustine, Luther had joined an important mendicant order, which by the middle of the 15th century had over 2,000 chapters. As a result of reforms carried through in 1473, the house at Erfurt, to which Luther went, accepted the strict, observant interpretation of the rule. Under Johann von Staupitz, Luther's mentor and vicar general to the order, a revised constitution was made in 1504. Luther made his profession as a monk in September 1506 and was then prepared for ordination. He was ordained priest in April 1507 and his first mass took place at the beginning of May. He had studied a treatise on the canon of the mass by a famous Tübingen Nominalist Gabriel Biel (d. 1495), who, like other "modern" Nominalists, claimed that only named particulars exist and that universal concepts are formed through intuition, and he approached the ceremony with awe. To this occasion his father came with a group of friends, and Luther took this first opportunity to explain personally the imperious nature of his vocation. His father's disgruntled retort, "Did you not read in Scripture that one shall honour one's father and mother?" struck deep into his memory.

Wittenberg University. Luther was selected for advanced theological studies; some of his university teachers were Nominalists of the "modern" way of the English philosopher theologian William of Ockham, whose views undercut the prevailing rationalism of Scholasticism, the school of thought founded in the 11th century in an attempt to reconcile revelation with reason. In 1508 Luther went to the University of Wittenberg (founded 1502), where, though Ockhamism had a foothold, the school of Realism that claimed that universals exist and can be known by reason was championed by scholars such as Martin Pollich. The little town was a contrast to Erfurt, but at least the university was young and forward-looking, and to its comparative remoteness Luther would one day owe his life. The Schlosskirche (called the Church of All Saints) was closely connected with the university, and the elector of Saxony, Frederick III the Wise (1463-1525), lavished generous patronage on both. In March 1509 Luther took the degree of *baccalaureus biblicus* at Wittenberg, returning to Erfurt for his next degree, of *sententiaris*, which involved expounding on the *Sentences*, a medieval theological textbook by Peter Lombard. He had begun his teaching with a course on Aristotle's *Nicomachean Ethics* and now began his career as a theologian with lectures on the *Sentences*. Some of his notes have survived, and if their theology is unexciting there is apparent an acid vehemence at the intrusion of philosophy and above all of Aristotle into the realm of theology.

Johann von Staupitz, vicar general of the German Augustinians, was very important in Luther's career as his teacher, friend, and patron. Staupitz seems to have been theologically trained as a Thomist (Realist) and was also influenced by the Augustinian tradition of his order, though his theology shows elements derived from the conflation in the late 15th century of the *devotio moderna* (modern devotion, a term used to describe the spirituality of the Brethren of the Common Life) with German mysticism.

The thunderstorm experience

Influence of Johann von Staupitz

His attempt to revive stricter discipline and to unite the observant and conventual Augustinians in Germany led to dispute, and Luther was one of two monks chosen to go to Rome to present the appeal of some dissident houses. He made the journey, the longest of his life, probably late in 1510, and his earnestness was shocked by the levity of the Roman clergy and by the worldliness so evident in high places. The appeal failed, and Luther returned to become a loyal supporter of Staupitz.

Staupitz became interested in his gifted pupil and, perhaps alarmed by his introspectiveness, encouraged him to proceed to his doctorate and to a consequent public teaching career. Luther took his D.Th. on October 19, 1512. The degree was important for Luther, with its implications of public responsibility. He soon took on the duties of a professor in succeeding Staupitz in the chair of biblical theology. This was his lifelong calling, and the exposition of the Bible to his students was a task that called forth his best gifts and energies, one that he sustained until ill health and old age made him relinquish it at the end of his life. In between lectures, in a manner of speaking, he began the Protestant Reformation.

Religious and theological questions. Meanwhile, Luther's own religious and theological difficulties were becoming acute. He had entered into the search for evangelical perfection with characteristic and serious zeal, and sought exactly to fulfill the rule of his order. Nonetheless, he soon found himself in problems difficult for him to understand, struggling against uncertainties and doubts, unhappily bearing a crippling burden of guilt, which neither the sacramental consolations (e.g., the Lord's Supper and penance) of the church nor the wise advice of skilled directors was able to assuage. This distress, which had its centre in his unquiet conscience, brought him into states of anxiety and despair. Nor were his difficulties lessened by the emphases of the Ockhamist theology, which encouraged an extroverted moralism, stressed the human will, and left aspects of uncertainty at the very points where Luther needed most to be reassured. "Temptation" (*Anfechtung*) was to become an important word for Luther's theology, a term that suggests the fight for faith, of which Staupitz could say that such experiences were meat and drink to Martin Luther. These inward, spiritual difficulties were enhanced by theological problems.

Discovery of "the righteousness of God." At the entrance to the world of the thought of St. Paul, Luther was halted—the road blocked by a word that intensified his difficulties to an almost intolerable degree. This was the conception of the "righteousness of God." His sombre childhood piety had made him intensely aware of God's judgment, and as a lecturer in the arts faculty at Wittenberg he had had to expound the Hellenic conception of justice, as he found it in the *Nicomachean Ethics* of Aristotle. Encouraged by the use of *justitia* ("righteousness" or "justice") in the works of several Nominalists, he came to think of God's justice as being primarily the active, punishing severity of God against sinners—i.e., in particular actions. It was for him a final aggravation of his trouble that in Rom. 1:17 it is asserted that the justice of God is revealed in the gospel. Thus, Luther concluded, the divine demand was shown as extending beyond outward obedience to the Law, revealed in the Commandments, to purity of heart, to inward motive and intention, so that grace itself became a demand and an exaction. Such a God could be feared but not loved, could be obeyed out of constraint but never with that happy spontaneity that Luther felt to be of the essence of Christian obedience.

Luther's inner conflict. To Luther's sense of failure to obey the Law was added the feeling of hypocrisy, which drove him to the edge of what moral theologians described as "open blasphemy." In 1545, in a celebrated autobiographical fragment that he prefaced to his complete works, he thus described his feelings:

For however irreproachably I lived as a monk, I felt myself in the presence of God to be a sinner with a most unquiet conscience, nor could I believe that I pleased him with my satisfactions. I did not love, indeed I hated this just God, if not with open blasphemy, at least with huge murmuring, for I was indignant against him, saying "as if it were really not

enough for God that miserable sinners should be eternally lost through original sin, and oppressed with all kind of calamities through the law of the ten commandments, but God must add sorrow on sorrow, and even by the gospel bring his wrath to bear." Thus I raged with a fierce and most agitated conscience, and yet I continued to knock away at Paul in this place, thirsting ardently to know what he really meant.

Thus, the dilemma. Illumination came at last, as in prayer and meditation he pondered the text, examining the connection of the words.

At last I began to understand the justice of God as that by which the just man lives by the gift of God, that is to say, by faith, and this sentence, "the justice of God is revealed in the Gospel," to be understood passively, that by which the merciful God justifies by faith, as it is written. "The just man shall live by faith." At this I felt myself to have been born again, and to have entered through open gates into paradise itself.

There has been great controversy about this inner conflict, but it seems certain that there was for Luther just such a crisis as he later described and that it was resolved in the manner he narrates. There has also been argument about the novelty of this discovery. There is in fact a profound difference between the Hellenic conception of distributive justice and the biblical doctrine of the righteousness of God as a divine, saving activity displayed in the field of history and of human experience, and Luther had penetrated deeply into the Pauline vocabulary at this point. The accuracy of Luther's memory about this and, indeed, his integrity have sometimes been impugned, but the verdict of a modern Catholic historian, Joseph Lortz, may stand: that if the discovery were not new, it was at any rate "new for Luther."

Salvation as grace. Had Luther not written this account, it would have been necessary to conjecture something like it to account for the new importance that he gave to justification by faith, a priority it retained in the new theological framework of Protestantism. This became for him the nerve of the gospel, that salvation is to be thought of primarily in terms of grace, and of a divine gift; that God's free, forgiving mercy is displayed in Jesus Christ; that the conscience, forgiven and cleansed, may be at peace, and that the soul, free from the burden of guilt, may serve God with a joyful, spontaneous, creative obedience. In his translation of the Bible Luther came to add "alone" after the word "faith" (*sola fide*) in the verse "For we hold that a man is justified by faith apart from works of law" (Rom. 3:28) because he felt it was demanded by the German language. The word alone or only was retained by the Reformers after him because it seemed to safeguard this important doctrine against such perversions as might seem to make salvation dependent on human achievement or a reward for human merit.

Evaluation of Luther's experience of justification. This experience ought not to be isolated, for Luther speaks of other problems of vocabulary (e.g., the conception of "repentance," *poenitentia*), and it cannot be assumed that this was for him a catastrophic personal experience such as befell St. Augustine, who had a mystical experience of God in the garden at Milan, or the 18th-century founder of Methodism John Wesley, who had a conversion experience at Aldersgate Street, London. About the date of the occurrence there has been much controversy. The publication of Luther's early lectures led naturally to the examination of these firstfruits of the young professor. Though an early view that it must have occurred during the period of Luther's first lectures on the Psalms (1513–15) has been damagingly criticized, Luther's use of the many-sided allegorism of the Middle Ages, which often found three or four levels of meaning in a single text, his concentration on the one historical meaning, and the Christ-centred core of theology of justification have led some scholars to believe that the illumination must have come to him before his lectures on the Letter to the Romans (1515–16).

Something depends on how the discovery itself is assessed: if it was a discovery that justification is a gift, that it is to be taken passively rather than actively, then (as the reference to Augustine's *De spiritu et littera*—"Concerning the Spirit and the Letter"—suggests) Luther was

Justification by faith

Luther's sense of hypocrisy

hardly moving beyond the Augustinian framework and it is probably from an early period. If, on the other hand, it was the more mature discovery of the relation of saving faith to the Word of God, then it must be placed later, perhaps in 1518–19. Many scholars now tend in this later direction, and they emphasize how Luther's thinking was stimulated and redirected by the urgent pressure of the church struggle that began in 1517.

Luther's
growth in
theological
thought

The net gain of this chronological discussion has been to demonstrate how important is the whole period of Luther's development from 1509 to 1521, and that his technical vocabulary and the categories of his theology were in movement throughout the whole of this period. Certainly his great courses of lectures on the Psalms (1513–15), on Romans (1515–16), Galatians (1516–17), and Hebrews (1517–18) reveal the growing richness and maturity of his thought.

Luther as preacher and administrator. Meanwhile, his other duties had accumulated. From 1511 he had been preaching in his monastery and in 1514 he became preacher in the parish church. This pulpit became the centre of a long and fruitful preaching ministry wherein Luther expounded the Scriptures profoundly and intelligibly for the common people and related them to the practical context of their lives. Within his order, he had become prior, and, in April 1515, district vicar over 11 other houses. Thus he became involved in a world of practical administration and of pastoral care that gave him valuable experience, standing him in good stead in later years when a large part of his vast correspondence would be concerned with the care of the German churches and the cure of needy souls.

The new University of Wittenberg found it must take sides in an academic crisis that faced the European universities of that day, the tension between an old and a new academic program. Before Luther's advent Martin Pollich, a leading professor at Wittenberg, had shown himself hospitable to Humanist influences, despite his preference for the older Thomism. Now Luther took the lead in inaugurating a new program, involving the displacement of Aristotle and the Scholastic theologians by a biblical humanism that turned to the direct study of the Bible, using as tools the revival of Greek and Hebrew and a renovated Latin and as a dogmatic norm the "old Fathers" (the early Church Fathers, or teachers) and especially St. Augustine. Such a program Luther planned with the help of his senior colleague, Karlstadt, and his young friend Philipp Melancthon. In February 1517 he penned a series of theses against the Scholastic theologians, which he offered to defend at other universities. Though this attempt to export the Wittenberg program met with no success he could write in May that the battle was won at least in Wittenberg—"our theology, and that of St. Augustine reign." But if his theses remained dormant, a very different fate awaited those that he wrote later in that same year. He could hardly have thought that these would fire a train that would explode the Western Christian world.

LUTHER AS REFORMER

The indulgence controversy. *The nature of indulgences.* The nature and scope of indulgences had been more and more defined during the later Middle Ages, but there was still an element of that dogmatic uncertainty that has been called a theological weakness of the age. Indulgences were the commutation for money of part of the temporal penalty due for sin, of the practical satisfaction that was a part of the sacrament of penance, which also required contrition on the part of the penitent and absolution from a priest. They were granted on papal authority and made available through accredited agents. At no time at all did they even imply that divine forgiveness could be bought or sold, or that they availed for those who were impenitent or unconfessed. But during the Middle Ages, as papal financial difficulties grew more complicated, they were resorted to so often that the financial house of Fugger of Augsburg had to superintend the sacred negotiations involved in them.

The way was open for further misunderstanding when in 1476 Pope Sixtus IV extended their authority to souls in

purgatory. The appeal to cupidity and fear, the pomp and circumstance with which these indulgences were attended, the often outrageous statements of some indulgence sellers were a matter of complaint. Luther himself had frequently preached against these abuses, for his patron, the elector Frederick, had amassed a great collection of relics in the castle church at Wittenberg, to which indulgences were attached. But the immediate cause of Luther's public protest was an indulgence that Frederick had prohibited from his lands, though it was available in nearby territory. This was a jubilee indulgence, offering special privileges, the ostensible purpose of which was the rebuilding of St. Peter's basilica in Rome. By a secret arrangement, half of the German proceeds were to go to the young Albert, archbishop of Mainz, who was deeply in debt owing to his rapid promotion to and payment for a number of high ecclesiastical offices.

The
immediate
cause
of the
indulgence
controversy

The Ninety-five Theses. Of this Luther knew nothing until some time afterward. For him, the provocation lay in the extravagant claims of an old, tried hand at this kind of thing, the Dominican salesman of indulgences Johann Tetzel. With these claims in mind, Luther drew up the Ninety-five Theses, "for the purpose of eliciting truth," and fastened them on the door of All Saints Church, Wittenberg, on October 31, 1517, the eve of All Saints' Day and of the great exposure of relics there. These were tentative opinions, about some of which Luther himself was not committed. They did not deny the papal prerogative in this matter, though by implication they criticized papal policy; still less did they attack such established teaching as the doctrine of purgatory. But they did stress the spiritual, inward character of the Christian religion, and the first thesis, which claimed that repentance involved the whole life of the Christian man, and the 62nd, that the true treasure of the church was the most holy gospel of the glory and the grace of God, showed the author's intention. The closing section attacked the false peace, that "security," which as a young lecturer Luther had so often attacked, of those who thought of divine grace as something cheaply acquired and who refused to recognize that to be a Christian involved embracing the cross and entering heaven through tribulation. Luther sent copies of the theses to the archbishop of Mainz and to his bishop. And here the invention of printing intervened. Copies were circulated far and wide, so that what might have been a mere local issue became a public controversy discussed in ever widening circles.

Reaction to the Ninety-five Theses. The archbishop of Mainz, alarmed and annoyed, forwarded the documents to Rome in December 1517, with the request that Luther be inhibited, at the same time reprimanding the indulgence sellers for their extravagance. At the time, it seemed to many that this was simply another squabble between the Dominicans and the Augustinians. Colour was given to this belief by the counter-theses prepared by a theologian, Konrad Wimpina, that Tetzel had defended before a Dominican audience at Frankfurt at the end of January 1518. When copies of these reached Wittenberg in March they were publicly burned by excited students. At Rome the pope merely instructed Gabriel della Volta, the vicar general of the Augustinians, to deal with the recalcitrant monk through the usual channels, in this case through Staupitz. Luther himself prepared a long Latin manuscript with explanations of his Ninety-five Theses, publication of which was held up until the autumn of 1518; it is a document of some theological importance, and shows how far from superficial Luther's original protest had been. Meanwhile, the chapter of the German Augustinians was held at Heidelberg, April 25, 1518. Luther was relieved of his extra duties as district vicar, in the circumstances a great relief and intended as such. He found great comfort in the support of his friends, and was himself in great form, winning over two young men, Martin Bucer, a Dominican, and Theodor Bibliander.

Counter-theses and student reaction

At this period Luther's theology was most especially a "theology of the cross"; *i.e.*, a theology that stressed the revelation of Christ on the cross. According to Luther, the "theology of the cross" seems foolishness to the wisdom of the world and is opposed to the natural theology of di-

"Theology of the cross"

vine power and majesty, which he attacked as a Scholastic "theology of glory." Important for him at this time was the inward religion preached by the 14th-century German mystic Johann Tauler and a short 14th-century mystical tract, the *Theologia Germanica*, that he himself edited and published (1516–18). In these months, therefore, he lay great stress on the need for the Christian to share the cross of Christ, in suffering and in temptation. Though these stresses were to recede into the background of Luther's developing theology, they were to remain important for the radical Reformation, for which the *Theologia Germanica* would be an important and seminal document.

Involvement of Johann Eck. During Luther's absence, and perhaps catastrophically, his senior colleague, Karlstadt, had taken action that was greatly to widen the scope and publicity of the controversy. The scholar Johann Eck (1486–1543) of Ingolstadt, a man of some learning, and with a zest for disputation, with whom Luther was already in friendly contact through a common friend, became involved in the controversy. He had written some observations on the Ninety-five Theses for his friend, the bishop of Eichstätt, and these manuscript observations, the so-called Obelisks, reached Wittenberg shortly before Luther went off to his chapter at Heidelberg; Luther himself replied with a few "Asterisks," but Karlstadt, concerned to defend the Wittenberg program, sprang into the fray with 379 theses, adding another 26 before publication. In some of these Eck was impugned. The Dominicans continued to press for Luther's impeachment, and proceedings against him for heresy began to move slowly in Rome. Luther himself did not improve matters by publishing a bold sermon on the power of excommunication that made it clear that here was not a man who would accept unquestioned whatever might be decided by the pope in terms of some undefined plenitude of power.

The Augsburg interview, 1518. *Luther before Cajetan.* A papal citation summoning Luther to Rome was sent to the cardinal Cajetan (1468–1534), a renowned Thomist, at Augsburg. But at this perilous moment politics fatefully intervened, and the period during which the Luther affair might have been swiftly, drastically disposed of without wider disaster to the church was eroded by considerations of policy. The elector Frederick, as one of the seven prince electors of the Holy Roman Empire, was most important to the pope, in view of the imminent choice of a new emperor, and the pope could not afford to antagonize him. The result was that Luther was bidden to a personal interview with Cajetan at Augsburg. He arrived there on October 7 with an imperial safe-conduct. The discussion had moved from indulgences to the discussion of the relation between faith and sacramental grace (the unmerited gifts of God in such acts as Baptism and the Lord's Supper), when an argument developed between the two theologians about the meaning of the "treasure" of merits that the papal definition of Sixtus IV said that Christ had acquired, and the incensed cardinal dismissed Luther from his presence, telling him to stay away unless he would unconditionally recant.

Luther's flight from Augsburg. While Luther waited uneasily, the Saxon councillors reported rumours that he would be taken in chains to Rome. Eventually, bundled through a postern by his friends, he fled the city. Now he wrote an appeal from the pope to a general (or ecumenical) council and a full defense of his actions to his prince. Cajetan, meanwhile, lost no time in denouncing Luther to Frederick, who was in something of a dilemma, though it counted much for Luther that he had the admiring friendship of the elector's secretary, the Humanist Georg Spalatin. At this time, too, the Wittenberg theological faculty addressed the prince on Luther's behalf, pointing out that the fate of the university and its reputation would be involved in Luther's disgrace. At one moment, it seemed that Luther might have to depart, perhaps for France or Bohemia. There then appeared Karl von Miltitz, a papal diplomat, who applied "stick and carrot" tactics to the elector, dangling before him at one moment threats against Luther and at the next the signal compliment of the golden rose, symbol of high papal honour and recognition. The diplomat promised more than he could possibly

perform, and after an interview with him at Altenburg, in January 1519, Luther sensed this and came to distrust him. A papal definition about indulgences, issued at Cajetan's request, seemed to show that Luther had indeed put his finger on some fatal ambiguities.

The Leipzig disputation, 1519. *Debate between Luther and Eck.* At Augsburg Luther had been in touch with Eck and arrangements were made for a public disputation at Leipzig in the summer. This was to be in the first place a debate between Eck and Karlstadt, though Luther was Eck's ultimate objective, but the hostility of George, duke of Saxony (the elector Frederick's first cousin), toward the Reformer raised difficulties about Luther's participation. Eventually it was arranged that Eck should debate with the two Wittenberg theologians in turn, in the castle of the Pleissenburg, Leipzig, at the end of July. There was a preliminary pamphlet skirmish. The issue between Eck and Karlstadt was the Augustinian doctrine of grace and free will, and Karlstadt wished to meddle neither with indulgences nor with papal authority. Among the preliminary matters, the origin of the papal power was raised and so Luther turned to a study of church history and Canon Law in the fateful weeks before the debate. A large contingent from Wittenberg attended, and in the presence of the theologians from both universities, Duke George and notables of church and state, the debate began. Eck showed some skill in manoeuvring Luther into a position in which he cast doubt on the authority of the great General Council of Constance (1414–18), and also defended some of the propositions of Jan Hus, a Bohemian Reformer who had been declared a heretic at Constance and burned to death at the stake. Leipzig was a part of Germany with a strong feeling against Bohemia, and the admission was received as damaging, giving ground for Eck's loud boast that the disputation had been his personal triumph. Luther, who had earlier said of the debate that it had not begun in God's name and would not end in his name, left Leipzig somewhat shaken and disturbed by Eck's verbal manoeuvring.

Luther's questioning of authority. Eck was able to go off to Rome with new prestige to give sharpness to the process of Luther's official condemnation. Luther had now to examine the further implications of his actions to date, in relation to the authority of the church, of councils, and of Scripture; his correspondence shows that he was reaching something like a crisis in his attitude to papal authority. There had been a small pamphlet war after the disputation that made it plain that there was strong support for Luther among the Humanists in Germany and Switzerland. Luther himself became involved in controversy with diverse theologians of Leipzig, and if he now wrote in the vernacular with increasing power and violence, his polemical writings reveal also his deep perceptions of the issues between himself and contemporary theology. Two Catholic universities, strongholds of tradition, Cologne and Louvain, next condemned Luther's teaching. But polemic was not Luther's main concern, and his *Sermon von den guten Werken* ("Sermon on Good Works"), issued in June 1520, is an important exposition of the ethical implications of justification by faith. As a tract it deserves to be associated with Luther's more famous tract on Christian liberty issued in the next months. On June 15, 1520, there appeared the papal bull (a decree issued under the papal seal) *Exsurge Domine* or "Lord, cast out," against 41 articles of Luther's teaching, followed by the burning of Luther's writings in Rome. Eck and the Humanist diplomat and cardinal Girolamo Aleandro (1480–1542) were entrusted with the task of taking the bull to the cities of Germany.

The Reformation treatises of 1520. Eck and Aleandro were alarmed to discover how swiftly German opinion had moved to Luther's side. In contrast to his treatment the year before, Eck had to seek refuge in Leipzig from physical violence. Aleandro did what he could in agitated correspondence to shock the Curia (papal administrative bureaucracy) into realizing the grave danger facing the church in Germany. Luther's friends, aware of how precarious his position was, sought to moderate his violence, but he now moved well beyond their horizon. In Luther's

Luther's
relation-
ship to
Jan Hus

Dispute
over
meaning
of merits

“Address to the Christian Nobility”

own opinion of himself, he was far too temperate in view of all the ecclesiastical hypocrisy and offenses. The result was the defiant tracts of the summer of 1520. The first, the real manifesto, was his *An den christlichen Adel deutscher Nation* (“Address to the Christian Nobility of the German Nation”), addressed to the rulers of Germany, princes, knights, cities, under the young emperor Charles V. It argued that in the crisis, when the spiritual arm had refused to take in hand the amendment of the church and the often expressed grievances of the German people against Rome (*i.e.*, the papacy), it was necessary for the secular arm to intervene and call a reforming council. The document was ill arranged and tailed off, but it found deep response among sections of the nation, and in the next months Luther was carried along with the tide of national resentment against Rome.

His second treatise, *De captivitate Babylonica ecclesiae praeludium* (“A Prelude Concerning the Babylonian Captivity of the Church”), intended for clergy and scholars, was an act of ecclesiastical revolution. It inevitably estranged many moderate Humanists, for it reduced to only three (Baptism, the Lord’s Supper, and penance) the seven sacraments of the church, denied mass and attacked transubstantiation (the doctrine that the substance of the bread and wine is changed into the body and blood of Christ in the sacrament of the Lord’s Supper), made vehement charges against papal authority, and asserted the supremacy of Holy Scripture and the rights of individual conscience. The third work, dedicated to the pope, was, as a still, small voice after the uproar, a minor classic of edification, *Von der Freiheit eines Christenmenschen* (“Of the Freedom of a Christian Man”), which made clear the ethical implications of justification by faith, and showed that his thought and his public actions were connected by a coherent theological core. On December 10, 1520, the students lit a bonfire before the Elster Gate in Wittenberg, and as they fed the works of the canonists to the flames Luther added the papal bull (*Exsurge Domine*) against himself with suitable imprecation—“Because you have corrupted God’s truth, may God destroy you in this fire.”

Excommunication

In January 1521 the pope issued the bull of formal excommunication (*Decet Romanum Pontificem*), though it was some months before the condemnation was received throughout Germany. Meanwhile, the imperial Diet was meeting at Worms, and there was a good deal of lobbying for and against Luther. In the end, Frederick the Wise obtained a promise from the emperor that Luther should not be condemned unheard and should be summoned to appear before the Diet. This enraged Aleandro, who asserted that the papal condemnation was sufficient and that the secular arm had only to carry out its orders. It also alarmed Luther’s friends, who did what they could to dissuade him. Luther was firm in his determination to go, and began the journey in April 1521, undeterred by the news, on the way, that the emperor had ordered his books to be burned. What was meant to be the safe custody of a heretic turned out to be something like a triumphal procession, and when Luther entered Worms on April 16 he was attended by a cavalcade of German knights and the streets were so thronged as to enrage his enemies.

The Diet and Edict of Worms, 1521. *Luther’s defiance of the Diet.* In the early evening of April 17, 1521, Luther appeared before the notables of church and state and faced the young emperor Charles V, whom he found cold and hostile. A pile of writings lay before him, but when he was formally asked whether he acknowledged them, his legal adviser insisted that the titles be read. In view of the gravity of recantation, Luther asked for time to think, a request that may have taken his enemies off guard. A day’s respite was granted, and the following afternoon, in a larger hall, and before an even more crowded assembly, Luther reappeared. This time he could not be prevented from making a long speech. He distinguished between his writings: for the works of edification he need not and ought not to recant, for the violence of his polemic he would apologize, but for the rest he could not recant; and, as he went on to explain why, the demand was brusquely made for a plain, simple answer. This he now gave in words of unyielding defiance. He would recant if

convinced of his error either by Scripture or by evident reason. Otherwise he could not go against his conscience, which was bound by the Word of God. Though evidence is now tilted against the authenticity of the famous conclusion, “Here I stand. I can do no other,” it at least registers the authentic note of Luther’s reply in a moment that captured the imagination of Europe. There was a moment of confusion with Eck and Luther shouting, and then the emperor cut short the proceedings. Luther strode through his thronging enemies to his friends, his arm raised in a gesture of relief and triumph.

There followed a diplomatic flurry. It was evident that Luther had powerful friends; there was some sabre rattling from the knights and the peasant emblem appeared in the streets. There is evidence to support Luther’s boast that had he wished he could have started such a game that the emperor’s life would not have been safe. The radical Reformer and social revolutionary, Thomas Müntzer, later asserted that had Luther recanted the angry knights would have killed him. At any rate, Luther was now given what he had long asked for in vain, something like a real hearing before reasonably impartial judges, while he was kindly handled by the archbishop of Trier. But he could not now make even minor concessions, and the discussions broke down on the fallibility of councils. He was formally dismissed and departed under his safe-conduct.

Despite his spectacular moral triumph, Luther’s enemies, nonetheless, achieved something important at this point when a rump Diet passed the Edict of Worms. It declared Luther to be an outlaw whose writings were proscribed. The edict was to shadow him and fetter his movements all his days. It meant also that his prince must, for a time at least, walk delicately and could not publicly support his protégé. The result was the pretended kidnapping of Luther who was lodged secretly in the romantic castle of the Wartburg, near Eisenach.

Luther at the Wartburg. In this aerie among the trees Luther remained until March 1522. Known as Junker Georg, or Knight George, he dressed as a layman, grew a beard, and put on weight. The lack of exercise and the unwontedly rich diet brought on physical distress, whereas his mind, flung back on itself after months of crisis, knew intense reaction in a period of acute depression of the kind that Luther ranked high among temptations. But he was far from idle. He finished a beautiful exposition of the Magnificat (the song of Mary, the mother of Christ, in the liturgy) and prepared an edition of sermons on the Epistles and Gospels at mass, which he thought was perhaps his best writing. Although away from books, he wrote his ablest controversial piece, *Rationis Latomianae pro Incendiariis Lovaniensis Scholae sophistis redditae Lutheriana confutatio* (“Refutation of the argument of Latomus”—who was a member of the theological faculty of the University of Louvain), containing a luminous exposition of justification. Most important of all, he began to translate the New Testament from the original Greek into German. He did not believe that such work should be left to one mind, and soon enlisted his colleagues, notably Melancthon, in the enterprise. But Luther’s was the controlling genius, and the resulting New Testament (published in September 1522), like the Old Testament, translated from the Hebrew, which followed later (1534), was a monumental work, which had deep and lasting influence on the language, life, and religion of the German people. He had now to deal with some of the practical implications of his revolt. Private masses, celibacy of clergy, religious vows were no theoretical questions, but were themselves entangled in a network of legal, financial, and liturgical affairs. He wrote about these things forthrightly, and Spalatin tried in vain to hold up their publication, for in Wittenberg there were growing difficulties, and the prince, the university, and the cathedral chapter were all, for various reasons, anxious to go slowly.

Commotion in Wittenberg, 1521–22. *Radical reform.* There was a lively section of the town and of the university, however, that was determined to force the pace, and there were violent scenes in the streets and churches early in October 1521. Yet Luther, on a secret visit to his friends early in December, was not alarmed, and it

“Here I stand. I can do no other.”

Influence of Luther’s translation of the Bible on the German language

was his influence that led the Augustinians to decide, in the new year, that those of them who wished might return to the world. Two radical leaders now appeared, the incorrigible troublemaker Karlstadt and Gabriel Zwilling, an ebullient spellbinder from the Augustinians. When Karlstadt announced his betrothal to a girl of 16, and at Christmas administered Communion in both kinds (bread and wine) while dressed as a layman, attacked images in a violent tract and in innumerable theses denounced vows and masses, and demanded a vernacular liturgy, it was evident that here was a program that in timing and method differed from Luther's. Moreover, its appeal to Scripture was legalistic and made matters of necessity things that for Luther lay within the option of Christian liberty. In the new year, the town council issued a notable and pioneering ordinance regulating religion, public morals, and poor relief, a document that owes much to Luther's teaching and perhaps something to the initiative of Karlstadt. At the end of 1521 confusion was increased by the arrival of the so-called Zwickau prophets, radicals on the run from the town of Zwickau, who spoke impressively of revelations given them through dreams and visions, claiming that the end of the world was near and that all priests should be killed. A flustered and outmanoeuvred Melancthon wrote urgently for advice to Luther, who sent wise and calm counsel.

Restoration of balanced reform. In the next months the situation worsened and in March 1522 Luther returned to Wittenberg, explaining the reason for his disobedience to instructions in a justly famous letter to his prince. Then, deliberately habited as an Augustinian monk once more, he took charge of his town pulpit and in a powerful series of sermons redressed the balance of reform. In these important utterances, the difference between Luther's conservatism and the radical pattern of reform is made plain. Luther deplored the use of violence, for the Word of God must be the agent of reform. He believed that revolt could not take place without destruction and the shedding of innocent blood; that the real idols are in the hearts of men and if their hearts are changed the images on church walls must fall into disuse. Moreover, the pace of reform must take into account the unconverted, weaker brethren. From this time onward Luther fought a war on two fronts, against the Catholics and against those whom he lumped together as *Schwärmer* ("fanatics"). One result of the Wittenberg crisis was to slow down the practical reforms, and though Luther introduced a reformed rite (*Formula Missae* or "Formula of the Mass," 1523) it was not until 1526 that he provided a vernacular liturgy (*Deutsche Messe*, or "German Mass"). Throughout Germany the evangelical movement continued to grow, and it was apparent that the Edict of Worms would not be everywhere enforced. A Diet at Nürnberg, 1522–23, refused to suppress the evangelical preachers and demanded a reforming, national council; though Catholic pressure was stronger in the following year, the Diet again pressed for a council and would consent only to the enforcement of the edict "as far as possible."

The Peasants' War. Activities of the radical Reformers. On his journeys to and from Worms Luther had been dismayed by the evident social and political unrest. In the next months he wrote open letters, warning the rulers of Saxony and the councils of such cities as Strassburg of the danger that the new radical teaching would provoke revolution. In 1523 he made his own views of secular government plain in an important treatise *Von weltlicher Obrigkeit* ("Of Earthly Government"), in which he firmly asserted the duty of a Christian prince and the place of secular government within God's ordinances for mankind; he distinguished between the two realms of spiritual and of temporal government, through which the one rule of God is administered, and stressed the duty of civil obedience and the sinfulness of rebellion against lawful authority.

In Saxony the radical teachers posed a problem for their untheological rulers. In Orlamünde, after having been rebuked at Wittenberg, Karlstadt had converted the community to his own brand of mystical quietism. Luther made a preaching tour of the area at the request of his prince, and was greeted with hostility and ridicule. Luther

himself denounced such social evils as usury, but in Eisenach the fiery preacher Jakob Strauss conducted a violent campaign against usury and tithes. Most formidable of all, in the little town of Allstedt, Thomas Müntzer, an unruly genius, combined his own ingenious liturgical reforms with a program of holy war. Himself a former "Martinian" (or follower of Martin Luther), he not only shared Karlstadt's enthusiasm for the mystics but added an explosive element (perhaps influenced by Hussite teaching) that gave point to Luther's worst fears. Müntzer threatened revolution and claimed that God would rid the world of its shame. Luther's warnings and events themselves forced the rulers to take action, and in the summer of 1524 Müntzer fled and Karlstadt was exiled. Müntzer wrote in a pamphlet that Luther was nothing more than a shameless monk, "whoring and drinking," and called him Dr. Liar. Karlstadt also wrote a series of tracts against his former comrades, denouncing, among other things, the corporeal presence in the Eucharist. Luther replied in a devastating and profound treatise, *Wider die himmlischen Propheten, von den Bildern und Sakrament* ("Against the Heavenly Prophets in the Matter of Images and Sacraments"). He claimed that the radical Reformers sought glory and honour, not the salvation of men's souls.

Luther's response to the Peasants' War. In the summer of 1524 the Peasants' War had broken out in the Black Forest area. Their program was variously motivated. Their demands were for concrete medieval liberties connected with the game and forest laws or with tithes. Some of them drew on Catholic teaching, others on the theology of Zwingli and of Luther, who had set an example of successful defiance of authority, had been no respecter of dignities, and whose teachings about Christian liberty and a priesthood in which all believers shared were plainer than his subtle distinctions between two kingdoms. Thus, both where he was understood and where he was misunderstood, Luther's influence in the Peasants' War has to be taken into account. Some of the moderate peasants included Luther among possible arbitrators. He himself published in May 1525 the *Ermahnung zum Frieden* ("Exhortation for Freedom"), an analysis of the "12 articles" of the Swabian peasants, sympathizing with just grievances, criticizing the princes, but repudiating the notion of a so-called Christian rebellion: "My dear friends, Christians are not so numerous that they can get together in a mob." Luther also claimed that the worldly kingdom cannot exist without inequality of persons.

In the spring of 1525 the Thuringian peasants rose, with Thomas Müntzer among their leaders, and at first seemed likely to carry all before them. Faced with imminent political chaos, Luther wrote a brutal, virulent broadsheet, *Wider die räuberischen und mörderischen Rotten der andern Bauern* ("Against the Murdering and Thieving Hordes of Peasants"). The writing was less violent than Müntzer's hysterical manifestos, but it was bad enough. It appeared, however, as an appendix to his moderate tract about the "12 articles." Moreover, words written at the height of the peasant success read very differently after their collapse at the Battle of Frankenhausen, May 15, 1525, and in the bloody reprisal that followed. It was typical of Luther that he refused to climb down, to regain lost popularity, and neither thereafter nor at any time can he be accused of subservience to rulers. As he had once refused to become the tool of the knights, so he had never "taken up" the peasant cause. But he confirmed many peasants in their preference for the radical ideology, which was soon to find more peaceful coherence in the Anabaptist movement.

Watershed year, 1525. Luther and Erasmus. In other ways, too, 1525 was a watershed in Luther's career. At the height of the Peasants' War in June 1525, "to spite the devil" he had married Katherina von Bora, a former nun. He certainly needed looking after, and she proved an admirable wife and a good businesswoman. His home meant a great deal to him and was an emblem for him of Christian vocation, so that he included domestic life among the three hierarchies (or "orders of creation") of Christian existence in this world, the other two being political and church life. In the same year there came his open break with the great Humanist Erasmus. The differences between

Luther's
dislike for
violence

Luther's
view of
rebellion

The two
realms
theory

Luther's
marriage

the two men had long been apparent, and Erasmus, who found in Luther the type of violent, dogmatic mendicant theologian he had always detested, liked what he saw of the Reformation less and less. Nonetheless, both men had a common band of admirers and friends and entered the arena with reluctance. Erasmus, in his *De libero arbitrio*, or "Concerning Free Will" (1524), attacked Luther's doctrine of the enslaved will and provoked a resounding reply in Luther's *De servo arbitrio*, or "Concerning the Bondage of the Will" (1525), a one-sided, violent treatise that, nevertheless, includes profundities still fruitfully debated. In that year, too, Frederick the Wise died. The two men had met only once, but Luther owed much to this prince. The new ruler, the elector John, and his successor John Frederick were Luther's devout supporters and with other princes, notably Philip, landgrave of Hesse, and Albert of Brandenburg, formed a coherent group in the imperial Diet.

The Diets of Speyer. The hostility of Charles V to the Reformers and his devotion to the Catholic faith never altered, but he had to take account of political exigencies, his quarrels with the Pope and with the king of France, and the need for support against the Turks. At the Diet of Speyer in 1526, the Edict of Worms was suspended, pending a national council; in the interval it was ruled that each prince must behave as he could answer to God and to the emperor. Luther stated that there was no fear or discipline any longer and that everyone did as he pleased. As a result, it was possible to plan the reorganization of the Saxon Church, and a visitation was carried out by jurists and theologians (1527–28). Some scholars have seen a tension between Melancthon's *Instruktion für die Visitatoren*, or "Instructions for the Visitation" (1528), and Luther's comments, which may reveal his distrust of secular intervention in spiritual affairs; and though he thoroughly approved of the development of the evangelical *Landeskirchen* ("territorial churches"), there were to be aspects of Lutheranism that blurred rather than reflected Luther's theological distinctions. At the second Diet of Speyer in 1529, renewed Catholic pressure led to the reversal of earlier concessions, drawing from the evangelical princes, and from a number of cities, a protest that won them, for the first time, the name Protestant.

The eucharistic controversy. *Doctrinal differences among the Reformers.* Doctrinal differences about the Eucharist broke the common evangelical front. Though all the Reformers repudiated the sacrifice of the mass, they were deeply divided about the nature of the divine Presence. Luther, with simple biblicism, insisted that Christ's words "This is my body" must be literally interpreted, because allegory is not to be used in interpreting Scripture unless the context plainly requires it. Karlstadt's fanciful argument (that the word this referred not to bread and wine but the Lord's physical body) was soon dropped. Zwingli won many to his view that "is" must be taken as "means," and his learned friend, the Humanist John Oecolampadius, brought support from the early Church Fathers for a spiritual Presence and stressed the idea of the 2nd-century Tertullian that "body" meant "sign of the body." Thus, the initial debate was about interpretive principles, about the words of institution, though the scriptural argument moved to the relevance or irrelevance of the Gospel According to John (e.g., "he who eats my flesh and drinks my blood has eternal life" [John 6:54]).

The debate turned to the intricate matter of Christology (i.e., doctrine of Christ). Zwingli insisted on the distinction between the two natures of Christ and that because it is the property of a human body to be in one place, Christ's human body was not here but in heaven. Luther, on the other hand, stressed the indivisible unity of the one Person of Jesus Christ, the mediator. Without going into a metaphysical doctrine of "ubiquity," or Presence everywhere (which was developed by other Lutherans), he asserted that Christ is present wherever he wills to be and that we are not to think of him in heaven "like a stork in a nest." Martin Bucer and the Strassburg theologians echoed the more positive stresses of the Swiss, and Bucer used the Realist language of the early Church Fathers to support a true, spiritual Presence. Luther's treatise *Dass*

diese Worte Christi "Das ist mein Leib" noch fest stehen wider die Schwärmgeister ("That these words of Christ 'This is my Body' still stand firm against the Fanatics," 1527) showed that in three years of controversy he had not budged. Zwingli's Latin tract *Amica exegesis* ("A Friendly Exegesis" 1527) was far less amicable than the title suggests and brought a great outburst from Luther, the impressive *Vom Abendmahl Christi, Bekenntnis* ("Confession of the Lord's Supper," 1528). This convinced Bucer that he had misunderstood Luther, who did not mean a local, confined Presence; and from then on he intensified his awkward, well-intended attempts to make peace.

The Marburg Colloquy and the Diet of Augsburg. The political advantages of a common front were obvious, not least to the vulnerable Zwingli and Philip, landgrave of Hesse, and the prince invited theologians of both sides to a private colloquy at Marburg in October 1529. Luther began by saying that in his opinion Zwingli did not know much about the gospel. When Zwingli asked if it was permissible for a Christian to ask how Christ could be present in the bread and wine of the Lord's Supper, Luther replied that if the Lord commanded him to eat crab apples and manure, he would do it because it was a command. After three days' debate, there was no agreement about the Eucharist, though the air had been cleared of many misunderstandings. But if the conference failed, there were agreements on other issues, and these might have been fruitful had not the coming imperial Diet caused the Wittenberg theologians to draw away from the Swiss. As an outlaw, Luther could not attend this fateful Diet of Augsburg and had to fidget in the castle of Coburg, leaving the care of the gospel to Melancthon, who did very well and produced in the Augsburg Confession (1530), one of the great documents of the Reformation as well as a normative confession of Lutheranism.

Luther used his influence to stiffen the elector against compromise, though from this time onward he could not refuse his consent to political Protestantism as it took a more and more military shape in the Schmalkaldic League, which was established by Protestant princes in preparation for armed resistance to Catholic aggression. The political situation again changed swiftly, however, and, confronted with the Turkish invasion, the Emperor agreed to a truce with the Protestants in the Religious Peace of Nürnberg (1532). This was a valuable breathing space, and its effects are evident in Luther's writings in the next years. Now, more and more, Luther left matters to the action of Melancthon. Opponents attempted to break up the friendship of the two. Luther said, regarding this matter, that if Melancthon would allow himself to be won over by their opponents, "he could easily become a cardinal and keep wife and child."

Growth of Lutheranism, 1530–46. *Melancthon's leadership.* Luther acquiesced in the eucharistic agreement—by which the south Germans reached agreement on the Lord's Supper—that the triumphant Bucer brought off with Melancthon in 1536 (the Wittenberg Concord), though Bucer was unable to widen the agreement and bring in the Swiss. When an English embassy from Henry VIII arrived to discuss joining the Schmalkaldic League, it was Melancthon who drew up the theological agenda (the Wittenberg Articles, 1535) with an ambiguous statement of justification of which Luther wrote, "this agrees well with our teaching." But he would not follow Melancthon when he thought he wrote too irenically about the papacy, and as the papal council loomed near he penned his own uncompromising Schmalkaldic Articles (1537).

Melancthon's great work in the field of education was to earn him the name preceptor of Germany, but Luther too was important in this matter. His open letter to the councillors of Germany about the need for schools (1524), and his published sermon *Dass man Kinder zur Schulen halten solle* ("On Keeping Children at School," 1530) show how wise and forward looking was his concern for education. He himself composed two important catechetical documents, the lovely classic, *Kleiner Katechismus* ("Small Catechism"), and *Grosser Katechismus* ("Large Catechism," 1529), for teachers and pastors.

In Wittenberg Luther had a group of able colleagues:

The
Augsburg
Confession

Luther
and
Zwingli

Justus Jonas, Johannes Bugenhagen, and Feliks Krzyzak (Cruciger). In scores of cities his disciples and friends spread the evangelical teaching that formed the Lutheran pattern of church life. Luther, though not pre-eminent as a liturgist, provided orders of worship from which numerous other *Kirchenordnungen* ("church orders") were derived. The influence of Luther's writings was everywhere felt in the Western Christian world. It was in Scandinavia that the Lutheran Church struck its deepest roots and won its most complete ascendancy, but it also had deep influence in Austrian and Hungarian lands. Luther realized the importance of hymns and encouraged his friends to write them. He wrote a score of fine hymns, four of which appeared in his first Protestant hymnbook in 1524. The famous "Ein feste Burg ist unser Gott" ("A Safe Stronghold Our God Is Still" or "A Mighty Fortress Is Our God") became almost an event in European history. During the last decade of his life, John Calvin (1509–64) was the rising portent in Switzerland, though Luther's personal contact with him was slight. He continued to attack bitterly the *Schwärmer* ("fanatics"), who then included besides the Anabaptists a number of radicals such as Kaspar Schwenckfeld, a Reformer who tried to mediate between various groups. Although he maintained to the end his view that error can be conquered only by the Word, Luther came to accept the punishment of the Anabaptists.

The affair of Philip of Hesse. In 1540 Bucer and Melancthon took the initiative in conniving at the deplorable bigamy of Philip of Hesse, but Luther was involved and had he willed could have stopped it. It would have been easy for Philip to remedy his incorrigible incontinence by taking a mistress, but this he refused to do, though his guilty conscience kept him from the sacrament. The desperate device, as a lesser of evils, was to grant him a secret dispensation to take a second wife. When the affair became public, Luther angrily threatened to expose the whole story. He himself was so far from lowering moral standards that in the next years he threatened to leave Wittenberg because public morals there were a shame on a city that had known the evangelical teaching so long. After a serious illness in 1537, he was an almost chronic invalid, prematurely aged, seldom free from discomfort, often in pain, and he brought his teaching career to an end with lectures on Genesis. In the last decade of his life, he had to witness the recovery of the papacy, which he thought to have been mortally wounded, in the preparations for the Council of Trent (1545–63), and the growing menace of Catholic military might. His last outstanding controversial treatise was *Von den Conciliis und Kirchen* ("Of Councils and Churches," 1539). Among his last writings, *Against the Anabaptists, Against the Jews, Against the Papacy at Rome, Founded by the Devil*, the most violent is the last, coarse and angry but still defiant.

Luther's last activities. Early in 1546 Luther was asked to go to Eisleben to mediate in a quarrel between two arrogant young princes, Counts Albrecht and Gebhard of Mansfeld. He was old and ill, but they were his *Obrigkeiten* ("authorities") to whom he owed obedience, and he set off in the snowy winter, leaving his wife stiff with anxiety. His letters to her teased her, comforted her, and spoke at last of a mission successfully accomplished. But he had overtaxed his strength, and in a few hours the chill of death came upon him. He died in Eisleben, where he was born, on February 18, 1546, and his body was interred in the Church of All Saints, Wittenberg. The great funeral orations by Bugenhagen and Melancthon, who knew him so well, are not simply panegyric. They witness that his intimates regarded him as a really great man, standing within the historic succession of prophets and doctors of the church, through whose life and witness the Word of God had gone forth, conquering and to conquer.

LUTHER AS THEOLOGIAN

Luther was no systematizer, like Melancthon or Calvin, though the dissensions among Lutheran theologians after his death, each appealing to one aspect of his thought, testify to the width, coherence, and delicate balance of Luther's own teaching. The basis of his theology was Holy Scripture; and, though the differences between his

own and Augustine's thought are important, Augustine must stand next to the Bible among the influences upon his mind. The doctrines of salvation were of prime importance for him, and here the two great, many-sided complex conceptions of the Word and of faith are important. His often subtle doctrine about civil obedience was not always understood by his later followers, and non-theological factors in German history perpetuated and, to a certain extent, even perverted this misunderstanding. His doctrine of Christian vocation in this world and the importance of human life in the world became part of the general Protestant and Puritan inheritance. In other matters—in the room allowed for Christian liberty, in his conception of the part played by law in Christian life, and in his insistence on the Real Presence in the Eucharist—his theology differs from the patterns that emerged in the Reformed (Presbyterian) churches, in Puritanism, and in the sects such as the Anabaptists.

MAJOR WORKS

In Latin

THEOLOGICAL WORKS: *Epistola Lutherana ad Leonem decimum summum pontificem. Dissertatio de libertate Christiana per autorem recognita* (1519; "Concerning Christian Liberty"); *De votis monasticis* (1521); *De captivitate Babylonica ecclesiae praeludium* (1520; "A Prelude Concerning the Babylonian Captivity of the Church"); *De servo arbitrio* (1525; "Concerning the Bondage of the Will").

CONTROVERSIAL WRITINGS: *B. Martini Lutheri theses Tezelio, indulgentiarum insitiori oppositas* (1517; Ninety-five Theses); *Rationis Latomianae pro incendiariis Lovaniensis scholae sophistis redditae Lutheriana confutatio* (1521).

EXEGESIS: *Enarrationes epistolarum et evangeliorum, quas postillas vocant* (1521).

In German

THEOLOGICAL WORKS: *Von den guten Wercken* (1520; "Of Good Works"); *Von weltlicher Ueberkeit, wie weyt man yhr gehorsam schuldig sey* (1523; "Of Earthly Government"); *Das diese wort Christi (Das ist mein leib etc.) noch fest stehen wider die Schwermgeyster* (1527; "That These Words of Christ 'This is My Body' Still Stand Firm Against the Fanatics"); *Vom Abendmal Christi, Bekenntnis* (1528; "Confession of the Lord's Supper"); *Von den Conciliis und Kirchen* (1539; "Of Councils and Churches").

CONTROVERSIAL WRITINGS: *An den christlichen Adel deutscher Nation* (1520; "Address to the Christian Nobility of the German Nation"); *Wider die hymelischen Propheten von den Bildern und Sacrament* (1525; "Against the Heavenly Prophets in the Matter of Images and Sacraments"); *An die Radsherrn aller Stedte deutsches Lands: Das sie Christliche Schulen aufsrichten und halten sollen* (1524); *Ermanunge zum Friden auff die zwelff Artikel der Bawrschaft ynn Schwaben* (1525); *Wider die mordischen wi reubischen Rotten der Bawren* (1525); *Wider Hans Worst* (1541); *Wider das Bapstum zu Rom vom Teuffel gestifti* (1545).

TRANSLATIONS AND EXEGESIS: *Das Neue Testament Deutzsch* (1522); *Biblia, das ist, die gantze Heilige Scrifft Deutzsch* (1534); *Das Magnificat verteuschet und ausgelegt* (1521).

OTHER WORKS (LITURGICAL): *Deutsche Messe* (1526). (DIDACTIC): *Der kleine Catechismus* (1559; "Small Catechism"); *Deutzsch Catechismus* (1529; "Large Catechism"). Among his hymns the most famous is probably "Ein feste Burg ist unser Gott" ("A Mighty Fortress Is Our God"). (E.G.R.)

BIBLIOGRAPHY

Luther's writings. Collections are the *Works of Martin Luther*, 6 vol., Philadelphia ed. (1915–32, reprinted 1982); and *Luther's Works*, American ed., edited by JAROSLAV PELIKAN and HELMUT T. LEHMANN, 55 vol. (1955–76), henceforth an indispensable tool for English study. In German the definitive edition is *D. Martin Luthers Werke: kritische Gesamtausgabe* (1883–), known as the Weimar edition. There is a single-volume anthology edited by JOHN DILLENBERGER, *Martin Luther: Selections from His Writings* (1961); also useful is E. GORDON RUPP and BENJAMIN DREWERY, *Martin Luther* (1970). The following are important volumes in the *Library of Christian Classics*: vol. 15, *Lectures on Romans*, ed. by WILHELM PAUCK (1961); vol. 16, *Early Theological Works*, ed. by JAMES ATKINSON (1962, reprinted 1980); vol. 17, *Luther and Erasmus*, ed. by E. GORDON RUPP and PHILIP S. WATSON (1969); and vol. 18, *Letters of Spiritual Counsel*, ed. by THEODORE G. TAPPERT (1955). Another important work is *A Commentary on St. Paul's Epistle to the Galatians*, ed. by PHILIP S. WATSON (1953).

Biographical and critical studies: PETER MANNS, *Martin Luther: An Illustrated Biography*, trans. from German (1982).

Luther's
hymns

Luther's
death

emphasizes the religious context. JOHN M. TODD, *Luther* (1982), is a popular biography. HEINRICH BORNKAMM, *Luther in Mid-Career, 1521-1530*, ed. by KARIN BORNKAMM (1983; originally published in German, 1979), examines Luther and his thoughts at midlife. MARK U. EDWARDS, JR., *Luther and the False Brethren* (1975), details the years between the Diet of Worms and Luther's death. H.G. HAILE, *Luther: An Experiment in Biography* (1980), concentrates on the last 10 years of his life. DAVID C. STEINMETZ, *Luther and Staupitz: An Essay in the Intellectual Origins of the Protestant Reformation* (1980), studies the influence on Luther of his early confessor and friend. ROLAND H. BAINTON, *Here I Stand!* (1950, reissued 1990), is a respected study. Also of interest are FRANZ LAU, *Luther* (1963; originally published in German, 1959); and W.J. KOOLMAN, *By Faith Alone* (1954; originally published in Dutch, 1946). PRESERVED SMITH, *The Life and Letters of Martin Luther* (1911, reprinted 1968), is the best of the older studies. A broad survey is E.G. SCHWIEBERT, *Luther and His Times* (1950). ROBERT HERNDON FIFE, *The Revolt of Martin Luther* (1957), portrays the young Luther. A brief account is E. GORDON RUPP, *Luther's Progress to the Diet of Worms, 1521* (1951, reissued 1964). WALTHER VON LOEWENICH, *Martin Luther: The Man and His Work* (1986; originally published in German, 1982), is an introductory analysis. GERHARD BRENDLER, *Martin Luther: Theology and Revolution* (1991; originally published in German, 1983), is a biography written from a Marxist perspective. A scholarly and readable interpretation of Luther is found in ERIC W. GRITSCH, *Martin—God's Court Jester: Luther in Retrospect* (1983). JAMES M. KITTELSON, *Luther the Reformer: The Story of the Man and His Career* (1986), makes Luther accessible to readers with little background in the history of the Reformation. BERNHARD LOHSE, *Martin Luther: An Introduction to His Life and Work* (1986; originally published in German, 1981), is also of special interest. The development of Luther, the man and the theologian, is assessed in HEIKO A. OBERMAN, *Luther: Man Between God and the Devil* (1989; originally published in German, 1982). MARTIN BRECHT, *Martin Luther*, 3 vol. (1985-93; originally published in German, 1983-87), is an in-depth portrait of the man and his times. Luther and his era are addressed in JAMES ATKINSON, *Martin Luther and the Birth of Protestantism*, rev. ed. (1982); A.G. DICKENS, *Reformation and Society in Sixteenth Century Europe* (1966, reprinted 1979); JOSEPH LORTZ, *The Reformation in Germany*, 2 vol. (1968; originally published in German, 1939); and WILHELM PAUCK, *Heritage of the Reformation*, rev. ed. (1961). MARK U. EDWARDS, JR., *Luther's Last Battles: Politics and Polemics, 1531-46* (1983), explores the influence of politics on Luther's thoughts, especially in his later years. Luther's politics are appraised in W.D.J. CARGILL THOMPSON, *The Political Thought of Martin Luther*, ed. by PHILIP BROADHEAD (1984). Critical studies on Luther's theology include GERHARD EBELING, *Luther: An Introduction to His Thought* (1970; originally published in German, 1964); PHILIP S. WATSON, *Let God Be God!* (1947, reissued 1970); E. GORDON RUPP, *The Righteousness of God* (1953, reissued 1963); HEINRICH BORNKAMM, *Luther's World of Thought* (1958; originally published in German, 1947); B.A. GERRISH, *Grace and Reason* (1962, reprinted 1979); REGIN PRENTER, *Spiritus Creator* (1953; originally published in Danish, 1944); and IAN D. KINGSTON SIGGINS, *Martin Luther's Doctrine of Christ* (1970). ALISTER E. MCGRATH, *Luther's Theology of the Cross: Martin Luther's Theological Breakthrough* (1985), focuses on the evolution of Luther's theology from 1509 to 1519. Luther's influence is traced in ERNST WALTER ZEEDEN, *The Legacy of Luther* (1954; originally published in German, 1950); and EDGAR M. CARLSON, *The Reinterpretation of Luther* (1948), a survey of Scandinavian Luther studies. Important studies written in languages other than English include KARL HOLL, *Gesammelte Aufsätze zur Kirchengeschichte*, 3 vol. (1921-28, reissued 1964); EMANUEL HIRSCH, *Lutherstudien*, 2 vol. (1954); RUDOLF HERMANN, *Gesammelte Studien zur Theologie Luthers und der Reformation* (1960); ERNST WOLF, *Peregrinatio*, 2nd ed., 2 vol. (1962); JOHANNES HECKEL, *Lex Charitatis*, 2nd ed. (1973); ERNST BIZER, *Fides ex auditu*, 3rd ed. (1966); OTTO HERMAN PESCH, *Die Theologie der Rechtfertigung bei Martin Luther und Thomas von Aquin* (1967, reprinted 1985); REINHARD SCHWARZ, *Fides, Spes, und Caritas beim Jungen Luther* (1962); and BERNHARD LOHSE, *Mönchtum und Reformation* (1963). Two psychological studies are PAUL J. REITER, *Martin Luthers Umwelt, Charakter und Psychose*, 2 vol. (1937-41); and ERIK H. ERIKSON, *Young Man Luther: A Study in Psychoanalysis and History* (1958, reissued 1993).

(Ed.)

Luxembourg

Strategic
location

The Grand Duchy of Luxembourg is a tiny sovereign state of only 998 square miles (2,586 square kilometres) that is bordered by Belgium on the west and north, France on the south, and Germany on the northeast and east. Luxembourg has come under the control of many states and ruling houses in its long history, but it has been a separate, if not always autonomous, political unit since the 10th century. The ancient Saxon name of its capital city, Lucilinburhuc ("Little Fortress"), symbolized its strategic position as "the Gibraltar of the north," astride a major military route linking Germanic and Frankish territories.

Luxembourg is a point of contact between the Germanic-

and Romance-language communities of Europe, and the grand duchy itself has three official languages: German, French, and Luxembourgian. The peoples of Luxembourg and their languages reflect the grand duchy's common interests and close historical relations with its neighbours. In the 20th century, Luxembourg became a founding member of several international economic organizations. Perhaps most importantly, the grand duchy was an original member of the Benelux Economic Union (1944), which linked its economic life with that of The Netherlands and of Belgium and would subsequently form the core of the European Community (EC).

This article is divided into the following sections:

Physical and human geography 314

- The land 314
 - Relief and soils
 - Climate
 - Settlement patterns
- The people 314
- The economy 315
 - Resources
 - Industry
 - Agriculture
 - International trade
 - Energy
 - Transportation

- Communications
- Administration and social conditions 316
 - Government
 - Health and welfare
 - Education
- Cultural life 317
- History 317
 - Ancient and medieval periods 317
 - Habsburg and French domination 317
 - Personal union with The Netherlands 318
 - Independent Luxembourg 318
- Bibliography 318

Physical and human geography

THE LAND

Relief and soils. The northern third of Luxembourg, known as the Oesling (Ösling), comprises a corner of the Ardennes Mountains, which lie mainly in southern Belgium. It is a plateau that averages 1,500 feet (450 metres) in elevation and is composed of schists and sandstones. This forested highland region is incised by the deep valleys of a river network organized around the Sûre (or Sauer) River, which runs eastward through north-central Luxembourg before joining the Moselle (or Mosel) River on the border with Germany. The Oesling's forested hills and valleys support the ruins of numerous castles, which are a major attraction for the region's many tourists. The fertility of the relatively thin mountain soils of the region was greatly improved with the introduction in the 1890s of a basic-slag fertilizer, which is obtained as a by-product of the grand duchy's steel industry.

The southern two-thirds of Luxembourg is known as the Bon Pays, or Gutland (French and German: "Good Land"). This region has a more varied topography and an average elevation of 800 feet. The Bon Pays is much more densely populated than the Oesling and contains the capital city, Luxembourg, as well as smaller industrial cities such as Esch-sur-Alzette. In the centre of the Bon Pays, the valley of the northward-flowing Alzette River forms an axis around which the nation's economic life is organized. Luxembourg city lies along the Alzette, which joins the Sûre farther north.

In the east-central part of the Bon Pays lies a great beech forest, the Müllerthal, as well as a sandstone area featuring an attractive ruiniform topography. The country's eastern border with Germany is formed (successively from north to south) by the Our, Sûre, and Moselle rivers. The slopes of the Moselle River valley, carved up in chalk and calcareous clay, are covered with vineyards and receive a substantial amount of sunshine, which has earned the area the name of "Little Riviera." Besides vineyards, the fertile soils of the Moselle and lower Sûre valleys also support rich pasturelands. Luxembourg's former iron mines are

located in the extreme southwest, along the duchy's border with France.

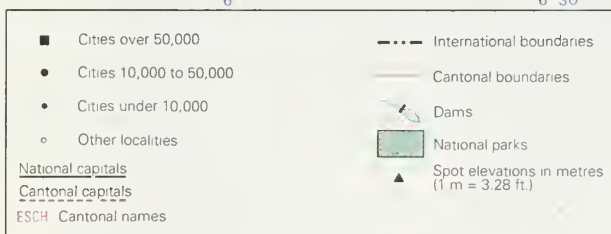
Climate. Luxembourg has a mild climate, with considerable precipitation. The north is slightly colder and more humid than the south. The mean temperatures in Luxembourg city range from 33° F (0.7° C) in January to 63° F (17° C) in July, but in the Oesling both extremes are slightly lower. The Oesling receives more precipitation than the Bon Pays, but the greatest amount, about 40 inches (1,000 millimetres), and the least, about 27 inches, fall in the southwest and southeast, respectively. The sheltered valley of the Moselle River benefits from a gentler and sunnier climate than does the rest of the duchy.

Settlement patterns. Northern Luxembourg is sparsely populated compared to the heavily urbanized and industrialized south. The north's rural population is clustered in villages of thick-set stone houses with slate roofs. The urban network in the south is dominated by the capital city, Luxembourg, which rises in tiers, with the upper (and older) section of the city separated from the lower-lying suburbs by the gorges of the Alzette and Petrusse rivers. A new quarter housing many European organizations nestles in a picturesque site carved into the river valley's sandstone cliffs. The second-largest city in Luxembourg, Esch-sur-Alzette, lies in the extreme southwest and is a traditional iron-making centre. Its growth, like that of the neighbouring iron and steel centres of Pétange, Differdange, and Dudelange, has slowed since the shrinkage of those industries in western Europe in the late 20th century. The remainder of the nation's population lives in towns and villages of relatively small size. Many of Luxembourg's villages date from ancient Celtic and Roman times or originated in Germanic and Frankish villages after about AD 400. In addition, many medieval castle villages continue to thrive, centuries after the castles themselves fell into ruin.

THE PEOPLE

Luxembourg has been one of the historic crossroads of Europe, and myriad peoples have left their bloodlines as well as their cultural imprints on the grand duchy. The Celts,

Main
regions



MAP INDEX

Political subdivisions

Capellen	49 38 N 6 00 E
Clerveaux	50 05 N 6 00 E
Diekirch	49 52 N 6 10 E
Echternach	49 46 N 6 22 E
Esch	49 30 N 6 02 E
Grevenmacher	49 40 N 6 20 E
Luxembourg	49 36 N 6 08 E
Mersch	49 45 N 6 06 E
Redange	49 48 N 5 54 E
Remich	49 32 N 6 20 E
Vianden	49 56 N 6 10 E
Wiltz	49 58 N 5 56 E

Cities and towns -

Bains
(Mondorf-les-Bains) 49 30 N 6 17 E
Bascharange 49 34 N 5 55 E
Bavigne, see Communauté du Lac de la Haute Sûre
Beckerich 49 44 N 5 53 E
Bertrange 49 37 N 6 03 E
Bettembourg 49 31 N 6 06 E
Bettendorf 49 52 N 6 13 E
Betzdorf 49 41 N 6 21 E
Bissen 49 47 N 6 05 E
Capellen 49 39 N 5 59 E
Clerveaux 50 03 N 6 02 E
Communauté du Lac de la Haute Sûre (Bavigne) 49 55 N 5 51 E
Consdorf 49 47 N 6 20 E

Dalheim 49 32 N 6 16 E
Diekirch 49 52 N 6 10 E
Differdange 49 31 N 5 53 E
Dudelange 49 28 N 6 06 E
Echternach 49 49 N 6 25 E
Erpeldange (Erpeldange-sur-Sûre) 49 52 N 6 07 E
Esch-sur-Alzette 49 30 N 5 59 E
Ettelbruck 49 51 N 6 07 E
Fouhren 49 55 N 6 12 E
Grevenmacher 49 41 N 6 27 E
Heiderscheid 49 53 N 5 59 E
Hesperange 49 34 N 6 09 E
Junglinster 49 43 N 6 15 E
Kayl 49 29 N 6 02 E
Kehlen 49 40 N 6 02 E
Kopstal 49 40 N 6 05 E
Larochette 49 46 N 6 14 E
Lintgen 49 43 N 6 08 E
Lorentzweiler 49 42 N 6 09 E
Luxembourg 49 36 N 6 08 E
Mamer 49 38 N 6 02 E
Mersch 49 45 N 6 06 E
Merttert 49 42 N 6 29 E
Mondercange 49 32 N 5 59 E
Niederanven 49 39 N 6 16 E
Pétange 49 33 N 5 53 E
Rambrouch 49 50 N 5 51 E
Redange (Redange-sur-Attert) 49 46 N 5 53 E
Remerschen 49 29 N 6 22 E
Remich 49 32 N 6 22 E
Rosport 49 48 N 6 30 E
Sanem 49 33 N 5 56 E
Schifflange 49 30 N 6 01 E
Steinfort 49 40 N 5 55 E
Strassen 49 37 N 6 04 E

Troisvierges	50 07 N 6 00 E
Useldange	49 46 N 5 59 E
Vianden	49 56 N 6 13 E
Walferdange	49 39 N 6 08 E
Wiltz	49 58 N 5 56 E
Wincrange	50 03 N 5 55 E
Wormeldange	49 37 N 6 25 E

Haute Sûre Lake reservoir	49 54 N 5 54 E
Haute Sûre Nature Park	49 52 N 5 50 E
Hou Forest	49 40 N 6 23 E
Mamer, river	49 45 N 6 06 E
Moselle (Mosel), river	49 32 N 6 22 E
Moselle Nature Park	49 33 N 6 20 E
Müllerthal Forest	49 48 N 6 18 E
Napoléonsgaard, peak	49 51 N 5 53 E
Oesling, region	49 50 N 6 00 E
Bon Pays (Gutland), region	49 45 N 6 10 E
Our, river	49 53 N 6 18 E
Our Nature Park	49 57 N 6 07 E
Schwaarzen Hiwel, peak	50 03 N 6 05 E
Sûre (Sauer), river	49 44 N 6 31 E
Syre, river	49 42 N 6 29 E

Physical features and points of interest

Alzette, river	49 52 N 6 07 E
Attert, river	49 49 N 6 05 E
Bassin Supérieur, reservoir	49 57 N 6 10 E
Bon Pays (Gutland), region	49 45 N 6 10 E
Buurgplaat, hill	50 10 N 6 02 E
Eisch, river	49 45 N 6 06 E
Ernz Blanche, river	49 52 N 6 16 E
Greng Forest	49 39 N 6 12 E

Moselle Nature Park	49 33 N 6 20 E
Müllerthal Forest	49 48 N 6 18 E
Napoléonsgaard, peak	49 51 N 5 53 E
Oesling, region	49 50 N 6 00 E
Our, river	49 53 N 6 18 E
Our Nature Park	49 57 N 6 07 E
Schwaarzen Hiwel, peak	50 03 N 6 05 E
Sûre (Sauer), river	49 44 N 6 31 E
Syre, river	49 42 N 6 29 E

the Belgic peoples known as the Treveri, the Ligurians and Romans from Italy, and especially the Franks were most influential. The dialect spoken by Luxembourg's native inhabitants is Luxembourgian, or Letzeburgesch, a Moselle-Franconian dialect of German that has been enriched by many French words and phrases. Most Luxembourgers speak French (used for most official purposes) and German (the lingua franca). There is a strong sense of national identity among Luxembourgers despite the prevalence of these foreign influences. Almost all of Luxembourg's native citizens are Roman Catholic, with a small number of Protestants, mainly Lutherans, and Jews.

National languages

Luxembourg has a higher proportion of foreigners living within its borders than does any other European country. This is chiefly the result of an extremely low birthrate among native Luxembourgers, which has led to a chronic labour shortage. Fully one-quarter of the total population is of foreign birth and consists mainly of Portuguese, Italians, and other southern Europeans, along with French, Belgians, and Germans. Among the foreign workers are many in the iron and steel industry, and numerous others work in foreign firms and international organizations located in the capital.

The 20th century has also witnessed a continual internal migration away from the countryside to urban areas, and the growth of Luxembourg's service sector at the expense of heavy industry has only accelerated this trend. Luxembourg city in particular continues to attract migrants from the rest of the country because of its vibrant banking and finance sector. The increasing concentration of the population in the southwest has led the government to try to locate some industries in rural areas. About two-thirds of Luxembourg's work force is engaged in trade, government, and other service occupations, while almost one-third of the work force is employed in industry and construction, and the rest in agriculture.

THE ECONOMY

Luxembourg's economy is notable for its close connections with the rest of Europe, since Luxembourg itself is too small to create a self-sustaining internal market. Luxembourg's prosperity was originally based on the iron and steel industry, which in the 1960s represented as much as 80 percent of the total value of exports. By the late 20th century, however, the nation's economic vigour stemmed chiefly from its involvement in international banking and financial services and in such noncommercial activities as hosting intra-European political activities. The result of this adaptability and cosmopolitanism is a very high standard of living; the Luxembourgers rank in the world second only to the Swiss in their standard of living and their per capita income.

Resources. Luxembourg's natural resources are quite modest. Its agriculture is not particularly prosperous, and its once-copious iron ore deposits had been exhausted by the 1980s. With the exception of water and timber, there are no energy resources. Indeed, Luxembourg has almost nothing that predisposes it to agricultural or industrial development. The roots of its economic growth lie in its use of capital and in the adaptability and ingenuity of its work force rather than in natural resources.



The town of Clervaux, in the Oesling, Lux.
E. Strechan/SuperStock

Industry. The production and export of iron and steel have long played major roles in Luxembourg's economy. Steel production was originally based on exploitation of the iron ore deposits extending from Lorraine into the southwestern corner of the grand duchy. This ore has a high phosphorus content, however, and it was not until the introduction of the basic Bessemer process in 1879 that the ore could be used for making steel. Thereafter Luxembourg's metallurgical industries grew and flourished. During the 1970s, however, the worldwide demand for steel slumped, causing the steel industry's portion of Luxembourg's gross domestic product to fall. In response to this crisis, the steel industry was restructured and merged into a single group called ARBED (Acieries Réunies de Burbach-Eich-Dudelange), and further measures aimed at increasing efficiency enabled Luxembourg's steelmakers to maintain their profitability. With the overall decline of steel production, however, Luxembourg's economy has become more dependent on factories owned by American-based and other multinational companies operating in the country. These factories primarily produce motor-vehicle tires, chemicals, and fabricated metals.

Luxembourg had become an international financial centre and a home to more than 160 banks by the late 20th century. It owes this position to a number of factors, perhaps chief of which is the government's own farsighted policies. In 1929 the government began to encourage the registering in Luxembourg of holding companies; these large corporations can control a number of subsidiary companies but are heavily taxed in many countries of the world. The liberal tax climate produced by the new policy led many industrial and financial corporations to maintain offices, often as their European headquarters, in Luxembourg city. The main offices of the European Investment Bank are there, as are the representatives of many banking institutions from around the world who keep in contact with the European Community (EC). Luxembourg city is also one of the capitals of the EC and as such is home to the European Court of Justice and several major EC administrative offices.

Agriculture. The agricultural resources of Luxembourg are quite modest. With the exception of livestock products, surpluses are scarce, and marginal soils in many parts of the country hinder abundant harvests. Most farming is mixed and includes both animal raising and gardening. Livestock and their by-products account for the bulk of agricultural production, with cattle raising having gained in importance at the expense of pig and sheep raising. Wheat, barley, and other cereal grains are the next most important products, followed by root vegetables. More

than three-quarters of the country's farms are smaller than 200 acres (50 hectares). The vineyards along the Moselle River produce some excellent wines.

International trade. Luxembourg's overall balance of payments is strongly positive, mostly because of the country's thriving financial sector. Most of the grand duchy's merchandise trade takes place with EC countries, and especially its three neighbours—Germany, Belgium, and France, which together receive about 60 percent of Luxembourg's exports and provide about 80 percent of its imports.

Energy. Luxembourg meets most of its energy needs with imports. Its only domestic source of power is the hydroelectricity obtained from several dams on its rivers.

Transportation. Luxembourg's internal road system is not extensive but is well maintained, and several highways link the country with its neighbours. A port at Mertert on the canalized Moselle River connects the grand duchy with the Rhine waterway system and provides it with an avenue for the international movement of goods.

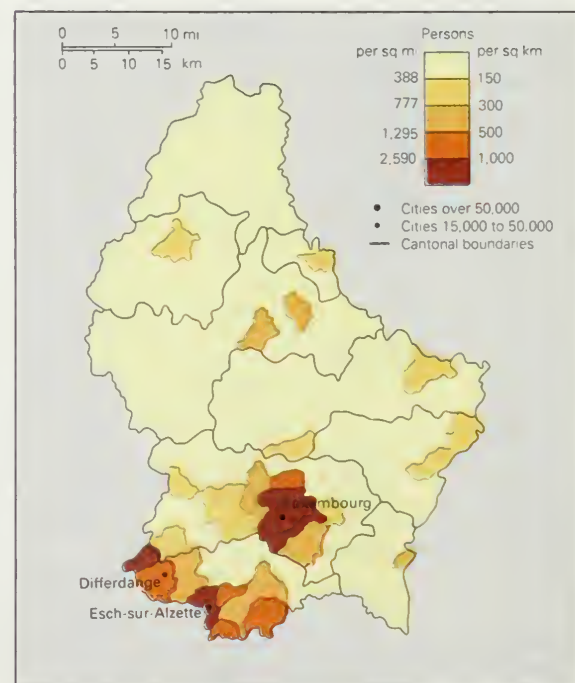
The government has operated the nation's railroads since World War II. They are modern, electrified, and mostly double-tracked. A major portion of international transportation to and from Luxembourg is by train, and the country is connected with its neighbours by a large number of lines. Findel Airport outside Luxembourg city has become a major European air terminal served by the lines of many countries. Luxair is the national airline.

Communications. Luxembourg's advanced telecommunications system provides it with close links both to EC countries and to other financial partners around the world, including Japan and the United States.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. The grand duchy is a constitutional monarchy with hereditary succession. Executive power authority lies with the grand duke, who appoints the prime minister. The powers of the grand duke are primarily formal, however. Actual executive power lies with the prime minister and his ministerial council, or Cabinet, who are responsible to the Chamber of Deputies. The members of this legislative assembly are elected by popular vote to five-year terms. Voting by all adult citizens, begun in 1919, is compulsory. Legislative elections have usually given rise to coalition governments formed alternatively by two of the three major parties: the Christian Social Party, the Socialist Workers' Party, and the Democratic Party. In addition, a Council of State named by the grand

Division of government responsibilities



Population density of Luxembourg.

duke functions as an advisory body. It is consulted on all draft legislation, advises the grand duke on administrative affairs, and serves as a supreme court in case of administrative disputes.

There are also three advisory bodies that are consulted before the passage of legislation affecting their particular area of the national life. The first of these consists of six confederations, three of which represent employers (commerce, guilds, and farmers) and three of which represent labour (workers, private employees, and civil servants). The second advisory group, the Social and Economic Council, has become a major committee for the examination of all projects. The third, the Immigration Council, advises the government on problems involving housing and the political rights of immigrants.

Justice is in the hands of magistrates appointed for life by the grand duke, the final appeal lying with the Superior Court of Justice. In the criminal court of assizes, six magistrates sit as jury as well as judge. Luxembourg is a member of the North Atlantic Treaty Organization (NATO) and has a small volunteer army. There is also a small paramilitary gendarmerie.

Luxembourg is divided administratively into three districts, each of which is headed by a commissioner appointed by the central government. Each district is in turn divided into cantons and subdivided into communes, or municipalities. Public works, health, and education are among the responsibilities of the communes, each of which is governed by an elected council and a mayor. These bodies also maintain liaison with the central government and act as its local agents.

Health and welfare. After World War I a broad system of social security and health services was introduced in Luxembourg to ensure maximum welfare protection to each citizen. Sickness benefits, in which patients pay only a small part of medical costs, as well as birth, family, and unemployment payments, are included in the plans. Housing conditions are generally comparable to those found in other western European countries. There has been some difficulty, however, in assimilating the many thousands of foreign workers and their families.

Education. Education is compulsory from age 6 to 15. The educational system offers a mix of primary and secondary schools run by state and local governments and by religious institutions. Considerable emphasis is laid on language studies. Instruction is initially given in German, and French is added in the second year. French completely replaces German in the classroom at the secondary level. There are no four-year universities in the grand duchy, so many young Luxembourgers obtain their higher education abroad. Luxembourg city does, however, have a campus that offers first-year university studies.

CULTURAL LIFE

The major cultural institution of Luxembourg is the Grand Ducal Institute, which has sections devoted to history, science, medicine, languages and folklore, arts and literature, and moral and political sciences. It functions as an active promoter of the arts, humanities, and general culture rather than as a conservator. The National Museum of History and Art has collections on the fine and industrial arts and on the history of Luxembourg. There is considerable public use of the National Library, the National Archives, and the Music Conservatory of the City of Luxembourg. The grand duchy also maintains cultural agreements with several European and other nations that provide it with the finest in the musical and theatrical arts. The Grand Orchestra of Radiotelevision Luxembourg is considered outstanding. There is an extensive market in Luxembourg city for works of painting and sculpture, both traditional and modern. The grand duchy's architectural heritage extends through practically the entire span of Europe's recorded history, from ancient Gallo-Roman villas to medieval castles, Gothic and Baroque churches, and contemporary buildings.

A small publishing industry exists, printing literary works in French, German, and Luxembourgian. The grand duchy's newspapers express diverse political points of view—conservative, liberal, socialist, and communist. Lux-

embourg's influence is felt far beyond its borders through the medium of Radiotelevision Luxembourg (RTL), a privately owned broadcasting company that transmits radio programs in five languages and television programs in two (French and German). RTL has a total European audience numbering as many as eight million persons. The government operates all postal and telegraph services in Luxembourg and has some control in the corporation that runs RTL.

For statistical data on the land and people of Luxembourg, see the *Britannica World Data* section in the BRITANNICA BOOK OF THE YEAR. (V.Bi./J.-P.E./J.M.G.)

History

ANCIENT AND MEDIEVAL PERIODS

The earliest human remains found in present-day Luxembourg date from about 5140 BC, but little is known about the people who first populated the area. Two Belgic tribes, the Treveri and Mediomatrici, inhabited the country from about 450 BC until the Roman conquest of 53 BC. The occupation of the country by the Franks in the 5th century AD marked the beginning of the Middle Ages in the locality. St. Willibrord played a very important role in the area's Christianization in the late 7th century. He founded the Benedictine abbey of Echternach, which became an important cultural centre for the region.

The area successively formed part of the Frankish kingdom of Austrasia, of the Holy Roman Empire under Charlemagne and Louis I the Pious, and then of the kingdom of Lotharingia. Luxembourg became an independent entity in 963, when Siegfried, Count de Ardennes, exchanged his lands for a small but strategically placed Roman castle lying along the Alzette River. This castle became the cradle of Luxembourg, whose name is itself derived from that of the castle, Lucilinburhuc ("Little Fortress"). Siegfried's successors enlarged their possessions by conquests, treaties, marriages, and inheritances. About 1060 Conrad, a descendant of Siegfried, became the first to take the title of count of Luxembourg. Conrad's great-granddaughter, Countess Ermesinde, was a notable ruler whose great-grandson, Henry IV, became Holy Roman emperor as Henry VII in 1308. This Luxembourg dynasty was continued on the imperial throne in the persons of Charles IV, Wenceslas, and Sigismund. In 1354 the emperor Charles IV made the county a duchy. In 1443 Elizabeth of Görlitz, duchess of Luxembourg and niece of the Holy Roman emperor Sigismund, was forced to cede the duchy to Philip III the Good, duke of Burgundy.

HABSBURG AND FRENCH DOMINATION

Along with the rest of the Burgundian inheritance, the duchy of Luxembourg passed to the Habsburgs in 1477. The division of the Habsburg territories in 1555–56 following Emperor Charles V's abdication put the duchy in the possession of the Spanish Habsburgs. In the revolt of the Low Countries against Philip II of Spain, Luxembourg took no part; it was to remain with what is now Belgium as part of the Spanish Netherlands. (For more specific information about the period, see THE NETHERLANDS.)

The duchy was able to remain aloof from the Thirty Years' War (1618–48) for a time, but in 1635, when France became involved, a period of disaster began in Luxembourg, which was wracked by war, famine, and epidemics. Moreover, the war did not end for Luxembourg with the Peace of Westphalia in 1648, but only with the Treaty of the Pyrenees in 1659. In 1679 France under Louis XIV began to conquer parts of the duchy, and in 1684 the conquest was completed with the capture of Luxembourg city. France restored Luxembourg to Spain in 1697, however, under the terms of the Treaties of Rijswijk. At the conclusion of the War of the Spanish Succession, by the treaties of Utrecht and Rastatt (1713–14), Luxembourg (along with Belgium) passed from the Spanish to the Austrian Habsburgs.

In 1795, six years after the beginning of the French Revolution, Luxembourg came under the rule of the French again. The old duchy was divided among three *départements*, the constitution of the Directory was imposed,

Radio-
television
Luxem-
bourg

Origin of
Luxem-
bourg

and a modern state bureaucracy was introduced. The Luxembourg peasantry were hostile toward the French government's anticlerical measures, however, and the introduction of compulsory military service in France in 1798 provoked a rebellion (the Kléppelkrieg) in Luxembourg that was brutally suppressed.

PERSONAL UNION WITH THE NETHERLANDS

French domination ended with the fall of Napoleon in 1814, and the Allied powers decided the future of Luxembourg at the Congress of Vienna in 1815. The Congress raised Luxembourg to the status of a grand duchy and gave it to William I, prince of Orange-Nassau and king of the Netherlands. William obtained a Luxembourg that was considerably diminished, since those of its districts lying east of the Our, Sûre, and Moselle rivers had been ceded to Prussia. The status of the grand duchy during this period was complex: Luxembourg had the legal position of an independent state and was united with The Netherlands only because it was a personal possession of William I. But Luxembourg was also included within the German Confederation, and a Prussian military garrison was housed in the capital city.

The rule of
William I

The standard of living of Luxembourg's citizens deteriorated during this period. Under Austrian rule, and especially from 1735 on, the duchy had experienced an economic expansion. From 1816-17 on, however, William I ignored the duchy's sovereignty, treating Luxembourg as a conquered country and subjecting it to heavy taxes. Consequently, it was not surprising that Luxembourg supported the Belgian revolution against William in 1830, and, in October of that year, the Belgian government announced that the grand duchy was a part of Belgium, while William still claimed the duchy as his own. In 1831 the Great Powers (France, Britain, Prussia, Russia, and Austria) decided that Luxembourg had to remain in William I's possession and form part of the German Confederation. Moreover, the Great Powers allotted the French-speaking part of the duchy to Belgium (in which it became a province called Luxembourg), while William I was allowed to retain the Luxembourgian-speaking part. Belgium accepted this arrangement, but William I rejected it, only to subsequently accede to the arrangement in 1839. From that year until 1867, the duchy was administered autonomously from The Netherlands.

INDEPENDENT LUXEMBOURG

William I negotiated a customs union for Luxembourg with Prussia, and his successor, William II, ratified this treaty in 1842. Against its own will, Luxembourg had thus entered into the Prussian-led Zollverein, or Customs Union, but the grand duchy soon realized the advantages of this economic union. Luxembourg subsequently developed from an agricultural country into an industrial one. Its road network was extended and improved, and two railway companies were begun that formed the basis for the national railway company founded in 1946.

Independent
Luxembourg

The restricted constitution that William II enacted for Luxembourg in 1841 did not meet the political expectations of its citizens. The Revolution of 1848 in Paris had its influence on the grand duchy, and William II that year enacted a new and more liberal constitution, which was in turn replaced by another constitution in 1856. In 1866 the German Confederation was dissolved, and Luxembourg became an entirely sovereign nation, though the Prussian garrison remained in the capital. Napoleon III of France then tried to purchase the grand duchy from William III. The two rulers had already agreed on the sum of five million florins when William III backed out because the Prussian chancellor, Otto von Bismarck, disapproved of the sale. The Great Powers soon came to a compromise (London, May 11, 1867): Prussia had to withdraw its garrison from the capital, the fort would be dismantled, and Luxembourg would become an independent nation. The grand duchy's perpetual neutrality was guaranteed by the Great Powers, and its sovereignty was vested in the house of Nassau.

On the death of William III of The Netherlands in 1890

without a male heir, the grand duchy passed to Adolf, duke of Nassau (d. 1905), who was succeeded by his son William (d. 1912). Neither Adolf nor William interfered much in Luxembourg's government, but William's daughter, the grand duchess Marie Adélaïde, was more assertive and eventually became highly unpopular with the people. In 1914 the neutrality of Luxembourg was violated by Germany, which occupied the grand duchy until the Armistice of 1918. During the war, Marie Adélaïde had tolerated the illegal German occupation, for which she was criticized by the Allied powers after the liberation. Marie Adélaïde was forced to abdicate in favour of her sister Charlotte in 1919. In a referendum a few months later, the public voted overwhelmingly against the establishment of a republic and in favour of retaining Charlotte as grand duchess.

In December 1918 the Allied powers had forced Luxembourg to put an end to its customs union with Germany. For the grand duchy this meant the loss of its best customer (for cast iron and steel) as well as its main supplier of coal. Luxembourg urgently needed a new economic partner, and, though the people preferred an economic union with France, the grand duchy was forced to negotiate with Belgium, since France declared itself uninterested in such a union. The Belgium-Luxembourg Economic Union (BLEU) was established in 1921 and provided for a customs and monetary union between the two countries. The economic climate in Luxembourg remained rather dreary during the interwar period though.

In May 1940 the German army invaded and occupied Luxembourg for the second time; however, this time the government refused to collaborate and, together with the grand duchess, went into exile. Luxembourg was placed under German rule, and the French language was banned.

After Luxembourg's liberation in September 1944, it took part in the new international organizations being formed by the victorious Allies, including the United Nations. Luxembourg also joined the new Benelux Economic Union (1944) formed between Belgium, The Netherlands, and itself. By taking part in the Brussels Treaty of 1948 and in the formation of NATO in 1949, Luxembourg abandoned its perpetual neutrality. The country improved its economic position by obtaining a sound position within the European Coal and Steel Community (1952) and within the European Economic Community (1957). Prince Jean, Charlotte's son, was installed as lieutenant-representant of Charlotte in 1961, and he inherited the throne in 1964 upon his mother's abdication. (V.L.)

Role in
European
unity

For later developments in the history of Luxembourg, see the BRITANNICA BOOK OF THE YEAR.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, sections 923, 961, 963, and the *Index*.

BIBLIOGRAPHY. An introduction to the country is provided by P. MARGUE *et al.*, *Luxembourg* (1984), in French, a well-illustrated work covering history, politics, ethnography, language and literature, natural history, and economy. Studies of the geography and economy of Luxembourg include J.M. GEHRING, *Le Luxembourg: un espace ouvert de l'Europe rhénane* (1977); PAUL WEBER, *Histoire de l'économie luxembourgeoise* (1950); and RAYMOND KIRSCH, *La Croissance de l'économie luxembourgeoise* (1971). Visual presentations of the land and of data on the country may be found in GUY SCHMIT and BERND WIESE, *Luxemburg in Karte und Luftbild* (1980), maps and aerial photos with text in German and French; and MINISTÈRE DE L'ÉDUCATION NATIONALE, *Atlas du Luxembourg* (1971). Works on Luxembourg's history are N. VAN WERVEKE, *Kulturgeschichte des Luxemburger Landes*, 2 vol. in 1 (1923-26); CHARLES J.P.A. HERCHEN, *History of the Grand Duchy of Luxembourg* (1950; originally published in French, 5th rev. ed. by N. MARGUE and J. MEYERS, 1947); E. DONCKEL, *Die Kirche in Luxemburg von den Anfängen bis zur Gegenwart* (1950); PAUL WEBER, *Histoire du Grand-Duché de Luxembourg*, 4th ed. (1961); *Manuel d'histoire luxembourgeoise*, 4 vol. (1973-77); JAMES NEWCOMER, *The Grand Duchy of Luxembourg: The Evolution of Nationhood, 963 A.D. to 1983* (1984); and GILBERT TRAUSSCH, *Le Luxembourg: émergence d'un état et d'une nation* (1989). For further resources on all aspects, see CARLO HURY and JULES CHRISTOPHORY (comps.), *Luxembourg* (1981), a bibliography. (J.M.G., V.L.)

Madagascar

The Republic of Madagascar (Malagasy: Repoblikan'i Madagasikara; French: République de Madagascar) lies off the southeastern coast of Africa. It occupies the fourth largest island in the world—after Greenland, New Guinea, and Borneo—with a surface area of 226,658 square miles (587,041 square kilometres). Located in the southwestern Indian Ocean, it is separated from the African coast by the 250-mile- (400-kilometre-) wide Mozambique Channel.

In spite of Madagascar's proximity to the continent, its population is primarily related not to African peoples but rather to those of Indonesia, more than 3,000 miles to the east. The Malagasy peoples, moreover, do not consider themselves to be Africans, but, because of the continuing bond with France that resulted from former colonial rule, the island has developed political, economic, and cultural links with the French-speaking countries of western Africa. French and Malagasy are the country's official languages. Madagascar remains a geographic and historical para-

dox, linked in practice to Africa but identified in feeling with Indonesia, which is so far away as to have hardly any awareness of Madagascar or to maintain any contemporary ties of substance with it. The animal life and vegetation of the island are equally anomalous, differing greatly from that of nearby Africa and being, in many respects, unique.

Although the coastlands have been known to Europeans for more than 400 years and to Arabs for much longer, recent historical development has been more intense and concentrated in the central plateau, which contains the capital city of Antananarivo (formerly Tananarive). The road network and communications are generally better on the plateau, where the majority of the inhabitants have received some school education and are professing Christians, while in the coastal areas the majority follow traditional religions and generally have not attended school. (A.S.)

The article is divided into the following sections:

Physical and human geography 319

The land 319

Relief

Drainage and soils

Climate

Plant and animal life

Settlement patterns

The people 322

Population groups

Demography

The economy 323

Resources

Agriculture, forestry, and fishing

Industry

Finance and trade

Administration of the economy

Transportation

Administration and social conditions 324

Government

Justice

The armed forces

Education

Health and welfare

Housing

Social and economic divisions

Cultural life 325

The cultural milieu

The state of the arts

Cultural institutions

The press and broadcasting

History 326

Early history 326

Madagascar before 1650

Political evolution 1650 to 1810

Early European contacts

The kingdom of Madagascar 327

Formation of the kingdom (1810–61)

Outside influences (1861–95)

The French period 327

The colonial period (1896–1945)

The French Union (1946–58)

The Malagasy Republic 328

Bibliography 328

Physical and human geography

THE LAND

Relief. Madagascar consists of three parallel longitudinal zones—the central plateau, the coastal strip in the east, and the zone of low plateaus and plains in the west.

Situated between 2,500 and 4,500 feet (800 and 1,400 metres) above sea level, the plateau has been uplifted and worn down several times and is tilted to the west. Three massifs are more than 8,500 feet high. The Tsaratanana region in the north is separated from the rest of the plateau by the Tsaratanana Massif, whose summit at an elevation of 9,436 feet (2,876 metres) is the highest point on the island. Ankaratra Massif in the centre is an enormous volcanic mass whose summit, Tsiafajavona, is 8,668 feet high. Ankaratra is a major watershed divide separating three main river basins. Farther south, Andringitra is a vast granite massif north of Tôlaïaro (Faradofay); it rises to 8,720 feet at Boby Peak.

The plateau slopes more regularly toward the extreme southern plain, but its boundaries to the east and west are more abrupt. To the east it descends in a sharp fault, by vertical steps of 1,000 to 2,000 feet. This cliff, which is called the Great Cliff or Cliff of Angavo, is often impassable and is itself bordered by the Betsimisaraka Escarpment, a second and lower cliff to the east, which overhangs the coastal plain.

Behind the scarp face are the remains of ancient lakes, including one called Alaotra. To the south the two steep gradients meet and form the Mahafaly and the Androy plateaus, which overhang the sea in precipitous cliffs. Toward the west the descent is made in a series of steps. In places, however, the central plateau is bordered by an impassable escarpment, such as the Cliff of Bongolava in the west-central part of the island. To the extreme north the plateau is bordered by the low belt of the Ambohitra Mountains, which include a series of volcanic craters.

The coastal strip has an average width of 30 miles. It is a narrow alluvial plain that terminates in a low coastline bordered with lagoons linked together by the Pangalanes (Ampangalana) Canal, which is some 400 miles long. To the south of Farafangana the coast becomes rocky, and in the southeast there occur many little bays. To the northeast is the deep Bay of Antongil (Antongila).

The western zone is between 60 and 125 miles wide. Its sedimentary layers slope toward the Mozambique Channel and produce a succession of hills. The inland (eastern) side of these steep hills dominates the hollows formed in the soft sediments of the interior, while the other side descends to the sea in rocky slopes. The coastline is straight, bordered by small dunes and fringed with mangroves. The currents in the Mozambique Channel have favoured the offshore deposit of alluvium and the growth of river deltas. On the northwestern coast there are a number of

The plateau

The coasts

estuaries and bays. This coast is bordered by coral reefs and volcanic islands, such as Nosy-Be (Nossi-Bé), which protects Ampasindava Bay. (Je.D.)

Drainage and soils. *Drainage.* The steep eastern face of the plateau is drained by numerous short, torrential rivers, which discharge either into the coastal lagoons or directly into the sea over waterfalls and rapids. They include the Mandrara, the Mananara, the Faraony, the Ivondro, and the Maningory. The more gently sloping western side of the plateau is crossed by longer and larger rivers, including the Onilahy, the Mangoky, the Tsiribihina, and the Betsiboka, which bring huge deposits of fertile alluvium down into the vast plains and many-channeled estuaries; the river mouths, while not completely blocked by this sediment, are studded with numerous sandbanks.

There are many lakes of volcanic origin on the island, such as Lake Itasy. Alaotra is the last surviving lake of the eastern slope. Lake Tsimanampetsotsa, near the coast south of Toliara (formerly Tuléar), is a large body of saline water that has no outlet.

Soils. The central plateau and the eastern coast are mainly composed of gneiss, granite, quartz, and other crystalline rock formations. The gneiss decomposes into red murrum, laterite, and deeper and more fertile red earths, giving Madagascar its colloquial name of the Great Red Island. Fertile alluvial soils in the valleys support intensive cultivation. There also are scattered volcanic intrusions that produce fertile but easily erodible soils. Lake Alaotra is a large sedimentary pocket in the central plateau containing some of the island's most productive farmland.

The western third of the island consists entirely of deposits of sedimentary rock, giving rise to soils of medium to low fertility.

Climate. The hot, wet season extends from November to April and the cooler, drier season from May to October. The climate is governed by the combined effects of the moisture-bearing southeast trade and northwest monsoon winds as they blow across the central plateau. The trade winds, which blow throughout the year, are strongest from May to October. The east coast is to the windward and receives a high annual rate of precipitation, reaching nearly 150 inches (3,800 millimetres) at Maroantsetra on the Bay of Antongil. As the winds cross the plateau, they lose much of their humidity, causing only drizzle and mists on the plateau itself and leaving the west in a dry rain shadow. The southwest in particular is almost desert, with the dryness aggravated by a cold offshore current.

The monsoon, bringing rain to the northwest coast of Madagascar and the plateau, is most noticeable during the hot, humid season. The wind blows obliquely onto the west coast, which receives a moderate amount of precipitation annually; the southwest, which is protected, remains arid. Annual rainfall drops from 83 inches on the northwestern island of Nosy-Be to 37 inches at Maintirano on the west coast to 14 inches at Toliara in the southwest. The plateau receives moderate rains, with 53 inches falling annually at Antananarivo and 48 inches at Fianarantsoa, which lies about 200 miles farther south.

July is the coolest month, with mean monthly temperatures around the island ranging from a low of about 50° F

The western rain shadow

The volcanic lakes

MAP INDEX

Political subdivisions

Antananarivo 19 00 s 47 00 E
Antsirafiana
(Antsiranana) 14 00 s 49 30 E
Fianarantsoa 21 30 s 47 00 E
Mahajanga 17 00 s 47 00 E
Toamasina 18 00 s 49 00 E
Toliara 21 00 s 45 00 E

Cities and towns

Ambalavao 21 50 s 46 56 E
Ambanja 13 41 s 48 27 E
Ambatolampy 19 23 s 47 25 E
Ambatondrazaka 17 50 s 48 25 E
Ambilobe 13 12 s 49 03 E
Ambohimahaso 21 07 s 47 13 E
Ambohitra 20 31 s 47 15 E
Ampanihy 24 42 s 44 45 E
Analavory 18 58 s 46 43 E
Andapa 14 39 s 49 39 E
Andoany
(Hell-Ville,
Hellville, or
Nosy-Be) 13 24 s 48 16 E
Ankarana
(Sosumav) 13 05 s 48 55 E
Antalaha 14 53 s 50 17 E
Antananarivo
(Tananarive) 18 55 s 47 31 E
Antsirabe 19 51 n 47 02 s
Antsirafiana
(Diégo-Suarez) 12 16 s 49 17 E
Antsohihy 14 52 s 47 59 E
Arivonimamo 19 01 s 47 11 E
Befandriana
Avaratra
(Befandriana-
Nord) 15 16 s 48 32 E
Besalampy 16 45 s 44 29 E
Betroka 23 16 s 46 05 E
Boriziny,
see Port-Bergé
Diégo-Suarez,
see Antsirafiana
Faradofay,
see Tôlaïfaro
Farafangana 22 49 s 47 50 E
Fénéry-Est,
see Fenoarivo
Antsinanana
Fenoarivo
Atsinanana
(Fénéry-Est) 17 22 s 49 25 E
Fianarantsoa 21 26 s 47 05 E
Fort-Dauphin,
see Tôlaïfaro

Hell-Ville,
see Andoany
Iharaña
(Vohemar or
Vohimarina) 13 21 s 50 00 E
Ihosa 22 24 s 46 07 E
Ivato 18 48 s 47 29 E
Maevatanana 16 57 s 46 50 E
Mahabo 20 23 s 44 40 E
Majunga,
see Mahajanga
Mahajanga
(Majunga) 15 43 s 46 19 E
Mahanoro 19 54 s 48 48 E
Maintirano 18 04 s 44 01 E
Mampikony 16 06 s 47 38 E
Manakara 22 08 s 48 01 E
Mananjary 21 13 s 48 20 E
Mandritsara 15 50 s 48 49 E
Maroantsetra 15 26 s 49 44 E
Marovoay 16 06 s 46 38 E
Miarinarivo 18 57 s 46 55 E
Moramanga 18 56 s 48 12 E
Morombe 21 44 s 43 21 E
Morondava 20 17 s 44 17 E
Nosy-Be,
see Andoany
Port-Bergé
(Boriziny) 15 33 s 47 40 E
Sambava 14 16 s 50 10 E
Sosumav,
see Ankarana
Tamatave,
see Toamasina
Tananarive,
see Antananarivo
Tangainony 22 42 s 47 45 E
Taolanaro,
see Tôlaïfaro
Toamasina
(Tamatave) 18 10 s 49 23 E
Tôlaïfaro
(Faradofay,
Fort-Dauphin or
Taolanaro) 25 02 s 47 00 E
Toliara (Toliary or
Tulear) 23 21 s 43 40 E
Tsiranomandidy 18 46 s 46 02 E
Vangaindrano 23 21 s 47 36 E
Vatomandry 19 20 s 48 59 E
Vohemar,
see Iharaña
Vohimarina,
see Iharaña
**Physical features
and points of interest**
Alaotra, Lake 17 30 s 48 30 E

Ambaro Bay 13 23 s 48 38 E
Ambatovaky
Special Reserve 16 50 s 48 40 E
Ambohitra
(Ambre)
Mountains 12 30 s 49 10 E
Ambohitra
(Ambre)
Mountains
National Park 12 37 s 49 09 E
Ambre, see
Bobaomby, Cape
Ampangalana
(Pangalanes)
Canal 22 48 s 47 50 E
Ampasindava
Bay 13 42 s 48 15 E
Andohahela
Nature Reserve 24 17 s 46 21 E
Andringitra
Massif 22 20 s 46 55 E
Andringitra
Nature Reserve 22 15 s 46 58 E
Androy Plateau 21 03 s 46 41 E
Angavo, Cliff of 18 30 s 48 00 E
Angavo
Escarpment 19 00 s 48 00 E
Ankarafantsika
Nature Reserve 16 09 s 47 02 E
Ankaratra Massif 19 25 s 47 12 E
Androntany (Saint
Sebastian), Cape 12 26 s 48 44 E
Antongila Bay of 15 45 s 49 50 E
Antsingin' i
Namoroka Nature
Reserve 16 28 s 45 20 E
Beampingaratra
Ridge 24 30 s 46 50 E
Bemaraha
Plateau 19 20 s 45 10 E
Betsiboka, river 16 03 s 46 36 E
Bobaomby
(Ambre), Cape 11 57 s 49 15 E
Boby Peak 22 11 s 46 53 E
Bongolava,
Cliff of 18 30 s 45 30 E
Faraony river 21 47 s 48 10 E
Fiherenana river 23 18 s 43 37 E
Ihotry, Lake 21 57 s 43 42 E
Ikopa, river 16 47 s 46 50 E
Indian Ocean 21 00 s 50 00 E
Isalo Massif 22 45 s 45 15 E
Isalo National
Park 22 45 s 45 10 E
Itasy Lake 19 04 s 46 47 E
Ivondro Massif 20 39 s 46 35 E
Ivondro, river 18 15 s 49 22 E
Kinkony Lake 16 08 s 45 50 E
Mahafaly Plateau 24 30 s 44 30 E
Mahajamba, river 15 33 s 47 08 E
Mahavavy, river 15 57 s 45 54 E
Makay Massif 21 15 s 45 15 E
Makira Plateau 15 30 s 49 15 E
Manambaho,
river 17 41 s 44 01 E
Mananantanana,
river 21 25 s 45 33 E
Mananara, river 23 21 s 47 42 E
Mananara, river 16 10 s 49 46 E
Mandrara, river 25 10 s 46 27 E
Mangoky, river 21 19 s 43 32 E
Mangoro, river 20 00 s 48 45 E
Mania, river 19 42 s 45 22 E
Maningory, river 17 13 s 49 28 E
Marojejy
(Marojejy) Nature
Reserve 14 30 s 49 44 E
Maromokotro
Peak 14 01 s 48 58 E
Marovoalavo
Plateau 16 45 s 48 48 E
Masoala
Peninsula 15 40 s 50 12 E
Matsiatra, river 21 25 s 45 33 E
Menarandra, river 25 17 s 44 30 E
Mitsio Island 12 54 s 48 36 E
Mozambique
Channel 20 00 s 41 00 E
Narinda Bay 14 55 s 47 30 E
Nosy-Be
(Nossi-Bé),
island 13 20 s 48 15 E
Onilahy, river 23 34 s 43 45 E
Pangalanes, see
Ampangalana
Canal
Sainte Marie
Island 16 50 s 49 55 E
Sofia, river 15 27 s 47 23 E
Tsaratana
Massif 14 00 s 49 00 E
Tsaratana
Nature Reserve 14 00 s 48 50 E
Tsiarafavona
Peak 19 21 s 47 15 E
Tsimanampetsotsa,
Lake, salt lake 24 07 s 43 45 E
Tsimanampetsotsa
Nature Reserve 24 08 s 43 47 E
Tsingin' i
Bemaraha Nature
Reserve 18 45 s 44 50 E
Tsiribihina, river 19 42 s 44 31 E
Vilanandro, Cape 16 11 s 44 27 E
Zahamena
Nature Reserve 17 40 s 48 50 E



Scale 1:7,692,000
 1 inch equals approx 121 miles
 0 25 50 75 100 125 mi
 0 50 100 150 200 km
 Lambert Conformal Conic Projection

<ul style="list-style-type: none"> ■ Cities over 150,000 ● Cities 50,000 to 150,000 • Cities under 50,000 ○ Other localities 	<p>National capitals</p> <p>Provincial capitals</p> <p>TOLIARA Provincial names</p> <p>— Provincial boundaries</p> <p>— Canals</p>	<ul style="list-style-type: none"> — Intermittent rivers — Reefs — Dams Swamps and marshes 	<ul style="list-style-type: none"> National parks ▲ Spot elevations in metres (1 m = 3.28 ft) <p>* Administered by France, claimed by Madagascar</p>
--	--	--	--

(10° C) to a high of about 78° F (26° C), and December is the hottest month, with temperatures between 61° and 84° F (16° and 29° C). Temperatures generally decrease with elevation, being highest on the northwest coast and lowest on the plateau.

Tropical cyclones are an important climatic feature. They form far out over the Indian Ocean, especially from December to March, and approach the eastern coast, bringing torrential rains and destructive floods.

Plant and animal life. Much of the island was once covered with evergreen and deciduous forest, but little now remains except on the eastern escarpment and in scattered pockets in the west. The plateau is particularly denuded and suffers seriously from erosion. The forest has been cut in order to clear rice fields, to obtain fuel and building materials, and to export valuable timber, such as ebony, rosewood, and sandalwood. About seven-eighths of the island is covered with prairie grasses and bamboo or small thin trees. There also are screw pines, palms, and reeds on the coasts. In the arid south of the island grow thorn trees, giant cacti, dwarf baobab trees, pachypodium succulents, and other xerophytes (drought-resistant plants) that are peculiar to the island.

Because of the island's isolation, many zoologically primitive primates have survived and evolved into unique forms. About 40 species of lemurs are indigenous to Madagascar. Several unique hedgehoglike insectivores have evolved, and there are many kinds and sizes of chameleons. Birds are numerous and include guinea fowls, partridges, pigeons, herons, ibis, flamingos, egrets, cuckoos, Asian robins, and several kinds of birds of prey. There are about 800 species of butterfly, many moths, and a variety of spiders. The only large or dangerous animals are the crocodiles, which occupy the rivers. The snakes, including the *do*, which is 10 to 13 feet in length, are harmless.

Inland waters contain tilapia (an edible perchlike fish), rainbow trout, and black bass. Marine fish and crustaceans abound on the coasts and in the lagoons, estuaries, and even in some upland streams. They include groupers, gilthead, tuna, sharks, sardines, whittings, crayfish, crabs, shrimps, mussels, and oysters. The coelacanth, referred to as a living fossil and once thought extinct for millions of years, inhabits offshore waters.

Settlement patterns. Despite the importance of intensive rice cultivation, the land is used primarily for pastoral purposes. Cattle are kept in all parts of the island. Fewer are found in the dense forest areas of the eastern escarpment, but elsewhere pastoralism predominates, most often coexisting with the cultivation of subsistence crops. On the plateau, the valley floors and irrigable slopes are mainly used for growing rice. The forest peoples traditionally grew hill rice, after cutting and burning the forest; this practice continues, although it is discouraged by the government, which promotes the establishment of permanent irrigated rice fields.

The older villages of the Merina and Betsileo were often perched on hilltops and defended by huge ditches. Today, villages have been rebuilt on lower ground, and hamlets and homesteads are scattered over the landscape. On the plateau, cattle enclosures are built of stone walls; the landscape is also dotted with funerary monuments, which take the form of beautifully carved wooden posts.

Small towns began to develop at the administrative centres of the island's several kingdoms at least by the 18th century. The most populous cities are Antananarivo, in the central plateau; Mahajanga (formerly Majunga), on the northwest coast; Fianarantsoa, in the southern plateau; Toamasina (formerly Tamatave), on the east coast; Antsiranana, in the north; Toliara, in the southwest; and Antsirabe, south of Antananarivo, mainly a tourist centre.

Racial or economic separation is not noticeable in the older cities such as Antananarivo and Fianarantsoa, but the newer towns are often divided into socioeconomic sectors. Because of recent internal migration, most of the cities are composed of a mixture of ethnic groups.

THE PEOPLE

Population groups. Madagascar has been inhabited by human beings for the relatively short period of only about

2,000 years. Language and culture point unequivocally to Indonesian origins, but there is no empirical evidence of how, why, or by what route the first settlers came to the island. Studies of the winds and currents of the Indian Ocean indicate that the voyage from Indonesia could have been made. It is assumed that the original peopling of the island, however sparse, was accomplished by a single cultural group, probably as the result of a single voyage.

There is also widespread evidence from linguistics, archaeology, and tradition of influence from Afro-Arab settlers on the coasts before AD 1000. There is slighter evidence of an Indian influence in vocabulary, but there is no trace of Hinduism in Malagasy culture and of orthodox Islām only in later coastal settlements.

The inhabitants of Madagascar speak Malagasy, which, written in the Latin alphabet, is a standardized version of Merina, an Austronesian language. Although there are numerous local variations of Malagasy, they are all mutually intelligible. French is also widely spoken and is officially recognized. It is used as a medium of instruction, especially in the upper grade levels, as is Malagasy.

The population is divided into about 20 ethnic groups, the largest and most dominant of which is the Merina people, who are scattered throughout the island. The name Merina (Imerina) is said to mean Elevated People, deriving from the fact that they lived on the plateau. The second largest group is the Betsimisaraka (The Inseparable Multitude), who live generally in the east. The third most numerous group is the Betsileo (The Invincible Multitude), who inhabit the plateau around Fianarantsoa. Other important peoples are the Tsimihety (Those Who Do Not Cut Their Hair); the Sakalava (People of the Long Valley); the Antandroy (People of the Thorn Bush); the Tanala (People of the Forest); the Antaimoro (People of the Banks); and the Bara (a name of uncertain origin). Smaller groups are the Antanosy (People of the Island); the Antaifasy (People of the Sand); the Sihanaka (People of the Lake); the Antakarana (People of the Rocks); the Betanimena (People of the Red Soil), who are now largely absorbed by the Merina; the Bezanozano (Those with Many Braided Hair); and the Mahafaly (the Joyful People). These ethnic names do not stand for clear-cut cultural boundaries, for in many cases one group fades imperceptibly into another. Moreover, the conventional translations are by no means reliable and most of the names themselves are of somewhat recent origin, probably crystallized and rigidified by the exigencies of colonial administration more than by the realities of indigenous culture. In no sense are these groups to be regarded as "tribes," a concept now considered invalid, nor are they composed of clans, but rather, in most cases, of endogamous and often non-unilinear descent groups.

About half of the population has been converted to Christianity, which is about equally divided between Protestantism and Roman Catholicism. A sizable community of Muslims also is found in the northwest. The rest of the people continue to practice their traditional religion, which is based upon ancestor worship. The dead are buried in tombs and are believed to reward or punish the living. There is a supreme being called Zanahary (the Creator) or Andriamanitra (the Fragrant One). There is also a belief in local spirits, and a complex system of taboos constrains Malagasy life.

Demography. More than 95 percent of the population is Malagasy. The major foreign communities are French, Comorian, Indian and Pakistani, and Chinese. Births greatly outnumber deaths, and the population is growing at a relatively rapid rate. Government policy, however, opposes any form of population control. One-half of the population, moreover, is under age 17, portending continued high growth rates well into the 21st century.

Emigration of the French, Comorians, Indians and Pakistanis, and Chinese in the late 20th century has significantly reduced their populations. There has been, however, no significant emigration of Malagasy peoples abroad.

The eastern part of the central plateau is the region of highest population density, and the eastern coastal plain has the second highest density. The eastern forest zone and the northeastern coast densities vary but rank as the

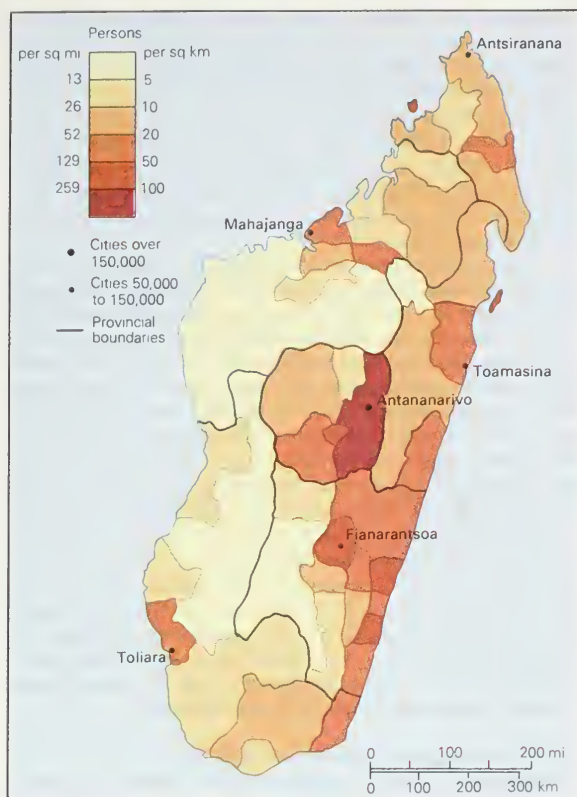
Indonesian origins

The grassland vegetation

Rice cultivation

Traditional religion

Population densities



Population density of Madagascar.

next most densely populated regions. Most of the western two-thirds of the country is sparsely inhabited.

The eastern half of the island contains almost all the major cities and towns. Antananarivo is the most populous city; perched on two precipitous mountain ridges, the old part of the city is dominated by the Manjakamiadana palace and has an extremely picturesque, almost medieval appearance.

THE ECONOMY

The control of Madagascar's economy by France, which had taken nearly half of the island's exports and supplied more than half of its imports, ended following the changes in government that occurred during the turbulent years 1972-75. Currency, banking, finance, loans, and economic planning had been influenced by accords with France and by French personnel in government, commerce, and technical assistance. After 1972-75 aid came from new sources, including the Soviet Union, the People's Republic of China, and North Korea, but most continued to come from France, other Western countries, and Japan.

Resources. The Great Red Island disappointed those who hoped to discover precious metals in large quantities. Considerable small-scale gold mining was conducted toward the end of the 19th century, by both French and Malagasy prospectors. There is a wide variety of semiprecious stones, but deposits are not of significant economic importance. Mineral deposits include chromite, which is found north of Antananarivo and in the southeast at Ranomena; ilmenite (titanium ore), found on the southeast coast at Tôlañaro, a source thought to represent the world's largest reserve of titanium; low-grade iron ore, found in scattered deposits in the southern half of the island; low-grade coal north of Toliara and inland from Besalampy; nickel near Fianarantsoa; and copper north of Ampanihy and near Ambilobe. There also are smaller deposits of zircon, monazite, bauxite, lead, graphite, quartzite, jasper, gold, uranothorianite, bentonite, kaolin, and alunite. The various gems include garnet, amethyst, tourmaline, beryl, and columbite.

The eastern evergreen forest and the scattered remnants of deciduous forest are valuable for their timber. Except for fish and crustaceans, the indigenous animal life is of little economic importance. Domestic animals include the humped zebu (cattle of Asiatic derivation, doubtless

brought from Africa), of great religious importance: goats; fat-tailed sheep; dogs; poultry; and pigs.

Although there are many narrow valleys and magnificent waterfalls, especially on the eastern escarpment, only a small number of them have been harnessed for electric power generation. Hydroelectric power stations provide about two-thirds of the country's electricity requirements; the remainder is supplied by coal-burning thermal stations. Many mines and factories also generate their own electricity with diesel- or steam-powered generators. Bituminous shales have been discovered at Bemolanga, oil at Tsimiroro, and natural gas off the coast of Morondava. The development of these resources remains commercially questionable, however.

Agriculture, forestry, and fishing. Rice occupies the largest share of the total crop acreage. Many varieties of dry, wet, and irrigated rice are grown in the central plateau; dry rice is also grown in the eastern forests and wet rice in the lower river valleys and along the estuaries, mainly by populations migrated from overpopulated parts of the plateau. Costly imports are still required.

Slash-and-burn techniques (the temporary clearance of land for agriculture) are used in the escarpment forest and along the east coast. In the river valleys of the west, cultivation is permanent; irrigation techniques are heavily utilized.

Sugarcane is grown on plantations in the northwest, around Mahajanga, and on the east coast near Toamasina. Cassava (manioc) is a staple grown all over the island, and potatoes and yams are cultivated mainly in the highland region of Ankaratra. Bananas are produced commercially on the east coast, and corn (maize) is grown mainly on the central plateau, in the south, and in the west. Fruits include apples, grapefruits, avocado pears, plums, grapes, oranges, litchis, pineapples, guavas, pawpaws, passion fruits, and bananas. Robusta coffee is grown on the east coast and arabica coffee on the plateau. Other significant crops are beans, peanuts (groundnuts), *pois du cap* (lima beans), coconuts, pepper, vanilla, cacao, sisal, raffia, tobacco, copra, cotton, and castor beans.

Cattle (mainly zebu) are distributed throughout the island. The large numbers of pigs, sheep, goats, chickens, ducks, geese, and turkeys are found mainly on the plateau. The hoarding of cattle as a sign of wealth and for religious sacrifice has frustrated government efforts to increase the use of cattle for domestic meat consumption and for export.

A significant area of the forest is degraded (*i.e.*, regenerated after repeated burnings, with many original species lost and smaller, fewer, less valuable species prevalent); the rest is wet or dry tropical forest. Major reforestation efforts have been undertaken, but, with about 80 percent of domestic fuel needs supplied by wood and charcoal, the country's total forested area continues to decline drastically.

Fisheries are poorly developed and depend mostly on small traditional fishing communities on the west coast; production is badly adjusted to the island's needs and to the potential market. In the late 20th century, however, the industrialization of marine fishing was begun, and rivers, lakes, and irrigated rice fields were stocked with breeding fish. The inshore waters are fished, while the distant offshore waters are neglected; inland distribution of the fish caught is poor, except in main towns. The bulk of the catch is composed of fresh fish and crustaceans; some fish are dried. There is also considerable raising of fish in the irrigated rice fields, mainly for home consumption.

Industry. Mining and quarrying are little developed. The most important products are graphite, mica, monazite, quartz, garnet, amethyst, salt, and chromite. Less valuable amounts of gold, tourmaline, citrine, beryl, ilmenite, columbite, zircon, and jasper are produced as well. Merina jewelers polish and set semiprecious stones at small workshops in most of the towns of the plateau.

The country's manufacturing consists mainly of rice mills and small industrial establishments. More than half of the manufacturing is located in and around Antananarivo. Products include wood, paper pulp, cotton cloth, fertilizer, oils, soap, sugar, cigarettes and tobacco, sisal

Agri-cultural produce

The fishing industry

Principal mineral deposits

rope and mats, bricks, processed foods, and beverages. A small printing industry and several motor-vehicle assembly plants are also located there.

Finance and trade. The Malagasy franc replaced the CFA (Communauté Financière Africaine) franc in 1963, and in 1973 Madagascar officially left the Franc Zone. The Central Bank issues all currency.

The value of exports is derived mainly from coffee, cloves and clove oil, vanilla, sisal, raffia, sugar, rice, pepper, tobacco, peanuts, and petroleum products. Despite their variety, many of these, like coffee, cloves, vanilla, and sisal, are threatened by world overproduction or by the manufacture of synthetic substitutes. Imports include crude petroleum, chemical and metal products, machinery, vehicles, electrical equipment, cotton textiles, and essential foods. After 1972, trade with France fell sharply, although France remains Madagascar's main trading partner. Other important trading partners include the United States, Japan, and Italy.

Administration of the economy. Before 1972 the government had established producers' cooperatives, which collected and processed most of the rice crop (at prices that were bitterly resented by the peasants); state farms intended to increase the commercial production of rice, cattle, coffee, oil palms, cotton, and silk; a rural development program; and a national consumers' cooperative with retail shops located in most towns.

After a period of prolonged domestic turbulence from 1972 to 1975, a new emergent military regime replaced the already constrained free-trade economy with one whose goal was to achieve "a socialist paradise under divine protection" by the year 2000. Nationalization made state corporations out of foreign firms, transforming the five French banks into three state banks for agriculture, industry, and trade. In addition, foreign insurance companies were converted into two state insurance corporations, and state monopolies were formed for import-export trading and shipping and for the textile, cotton, and power industries, as well as for the new agencies created for the extension of irrigation.

The economy responded with a relentless decline. Exports fell, inflation rose, and debt expanded from \$89 million in 1970 to \$250 million in 1978 and \$3.3 billion in 1988. Debt services took more than half of the country's export earnings, and imports nearly ceased for lack of foreign exchange. Industrial production fell to a quarter of the 1975 level, and foreign investment declined to almost nothing. Devaluation reduced the value of the Malagasy franc by 80 percent. Under these adverse conditions the country's infrastructure and social services deteriorated greatly. The rural population was reduced to subsistence levels, bartering with cattle and bags of paddy. By 1982 the country was technically bankrupt, forced to adopt a program of structural adjustment imposed by the International Monetary Fund and face a humiliating turnabout to reliberalize the economy. Abolition of the state monopolies was accepted, and any state enterprise that could not pay its way was threatened with privatization. Private business took years to regain any confidence, however.

Taxation is mainly indirect and is derived largely from various customs, import and export duties, and excise taxes. Direct taxes take the form of taxes on company income, registration fees, stamp duties, and personal taxes; the latter bear heavily on the peasants.

Most union members hold office jobs or work in industry. A popular and extensive trade union is the Union des Syndicats Autonomes de Madagascar, but the Sendika Kristianina Malagasy (Christian Confederation of Malagasy Trade Unions) claims more affiliated unions. The main employers' association is the Union des Syndicats d'Intérêt Économique de Madagascar.

Transportation. Transport facilities serve primarily the plateau and the east coast. Facilities are rudimentary on the western half of the island, although the country's best natural harbours are located there.

The majority of roads are unpaved. Roads down the eastern escarpment and across the western coastal strip, as well as minor roads everywhere, become impassable during the wet season. The main paved road runs south from

Antananarivo to Fianarantsoa, where it branches southwest to Toliara, southeast to Tôlanaro via Ihoay, and east to Mananjary and Manakara. Paved roads run east from Antananarivo toward Toamasina, west to Analavory, and north to Mahajanga and Antsiranana.

Two railways connect the plateau with the east coast; they run from Antananarivo to Toamasina and from Fianarantsoa to Manakara. Plateau routes run from Antananarivo south to Antsirabe, to Fianarantsoa, and north to Ambatondrazaka.

The main port is Toamasina, which has a fine deepwater harbour equipped with quay berths and directly linked to Antananarivo by rail, air service, and road. Mahajanga is second in importance but is accessible only to small ships of shallow draft; it has considerable dhow traffic with the Comoros. Antsiranana has one of the finest natural harbours in the world but is as yet too remote from the main centres of economic activity; it contains the former French naval base, arsenal, and dry dock and also has a small commercial port. The coastal lagoons and swamps of the Pangalanes Canal on the east coast, which are linked by artificial channels where necessary, provide a waterway that is 434 miles (698 kilometres) long.

Airfields are found throughout the island. The main international airport is at Ivato, near Antananarivo, and some international flights make secondary landings at Toamasina, Nosy-Be, and Mahajanga. Air Madagascar provides internal service; its international flights are supplemented by those of several foreign airlines.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. The Malagasy Republic became independent in 1960, having been an autonomous republic in the French Community from 1958 (see below *History*). The country was ruled from 1960 to 1972 through a bicameral legislature consisting of a national assembly, directly elected by universal adult suffrage, and a senate, one-third of which was nominated by the president and the rest of which was indirectly elected.

Between 1972 and 1975 Madagascar was under military rule. Radical political and economic reorganization was instituted in 1975, and a new constitution was implemented later that year for the renamed Democratic Republic of Madagascar. The constitution established a presidency; the Supreme Revolutionary Council, a body of members including the president, prime minister, military chiefs, some ministers, and other nominated leaders; a Cabinet, composed of the prime minister and other ministers; a unicameral legislature, the People's National Assembly, composed of deputies directly elected for five years; a Military Development Committee for advice on defense and economic policy; and a National Front for the Defense of the Revolution (Front National pour la Défense de la Révolution), grouping all parties willing to support the government and led by the Advance Guard of the Malagasy Revolution (Avant-garde de la Révolution Malgache) with its "Red Book" ("Boky Mena") manifesto. The government was fairly well-balanced ethnically, and the old polarization of coastal and plateau peoples was diminished, although the Merina continued to dominate the civil service, professions, and education.

There was a many-layered local electoral system in which thousands of local village communities elected delegates to associations of village councils, which then sent delegates to district groupings, which, in turn, sent delegates to six provincial councils. Power was in the hands of party delegates at all levels, under military oversight, rather than with the former bureaucratic hierarchy that existed before the reorganization.

With the subsequent failure of the government's utopian policies and the devastation and virtual bankruptcy of the economy, the constitution was relaxed, again allowing for opposition parties, of which two emerged very powerful. Pressing for abolition of the Supreme Revolutionary Council and the Military Development Committee, the opposition has also sought another new constitution.

Justice. The former Merina state that ruled the island throughout the 19th century had an elaborate system of laws, courts, and justice. The present Malagasy legal sys-

Rail
transportation

Economic
planning

"Red
Book"
manifesto

Customary law

tem, however, is based upon French codes and practices, and most judges and magistrates have had French training. There is a Constitutional High Court, an appellate court, criminal tribunals, and tribunals of first instance; there is also a criminal code, a code of criminal procedure, and a code of civil procedure. The customary law of the Merina and other ethnic groups is taken into account by state magistrates when judging marriage, family, land, and inheritance cases.

The armed forces. The army played no direct political role in Madagascar until 1972, although the presence of French army units had bolstered the former government. The French troops were withdrawn in 1973, and the French naval base at Antsiranana was handed over to Madagascar in 1975. The Malagasy armed forces consist of an army, a navy, and an air force. There is also a large paramilitary force, as well as a secret police. Detachments of local police are stationed at the headquarters of each administrative division, as well as in Antananarivo. The whole force is under unified command and falls within the responsibility of the ministry of the interior.

Education. The educational system consists of primary and secondary schools, technical institutes, teacher-training colleges, and a university system. Enrollment at the University of Madagascar (founded in 1955) and its five regional branches has increased dramatically in the late 20th century. There has been an increased use of the Malagasy language in teaching, although some coastal peoples have objected because of the language's close relationship to the Merina people.

The role of mission schools

The level both of school attendance and of educational attainment is higher on the plateau than in the coastal areas. Protestant and Roman Catholic missions have been providing education since the 19th century, and the missions continue to educate a large proportion of the schoolchildren, although the government now maintains official schools at all levels and is attempting to phase out private education. In the main towns there are other privately run schools, catering to those unable to enter either government or mission schools. About two-thirds of Madagascar's population is literate, and the majority of illiterate persons are female.

Health and welfare. Malagasy doctors began to practice Western medicine in 1880; and a medical school was established in Antananarivo. The health system includes principal and secondary hospitals, dispensaries, and medical centres. Medical personnel include doctors, as well as pharmacists, dentists, midwives, social assistants, visiting nurses, and health assistants.

Hospitals and specialists are mainly in the towns, apart from some rural hospitals run by Christian missions. Health insurance and other social benefits are available mainly to better-paid workers and professionals among the employed population.

The extension of health services is largely credited for the steady population increase. Infant mortality remains high, but infant deaths from malaria, which is endemic all over the island, have been cut by half. Debilitating parasitic diseases, such as schistosomiasis, an infection of the bladder or intestines, remain serious and are hard to control since their breeding grounds are the irrigated rice fields and the streams that feed them. Venereal disease is also widespread, especially in its incipient form.

Housing. Houses are typically rectangular and crowned with steeply angled roofs. In the rural areas, most houses are made of either mud and wattle or woven matting supported by poles. In the eastern forest, they are built of interlaced split bamboo and are thatched with palm, while, in the south, overlapping upright wooden planks are used for the walls. In the plateau, rural housing is constructed of earth blocks and thatched roofing, while upper-income and most urban housing consists of two- or three-story homes—typically with kitchen at the top, living quarters in the middle, and storage below—all surrounded by wide balconies supported by brick columns and crowned with steep tiled roofs. This is the lofty Indonesian style of architecture, transformed by new techniques contributed by the missionaries. The original style survives in the house of Andrianampoinimerina (reigned 1787–1810) at Ambo-

Indonesian style of architecture

himanga and reaches its apotheosis in the queen's palace built by Jean Laborde in the 19th century.

The government-sponsored housing authority conducts research into design, materials, and production methods and is seeking to promote inexpensive urban housing, but the problem of overcrowding is expected to increase with continued urban growth. The existence of a well-established craft of house construction, however, may successfully alleviate housing pressures without resorting to imported materials or relying on foreign enterprise.

Social and economic divisions. Traditionally, society was divided into three castes—the nobles, the freemen, and the former slaves and their descendants. These social distinctions are no longer strict and are manifest only on ceremonial occasions, such as weddings and funerals. They do, however, form the basis of other economic and social distinctions. During the 19th century, the Merina elite conquered the island, established themselves as rulers, and adopted Protestant Christianity; in the late 1800s, some became Roman Catholics. Under French rule in the 20th century, the Merina retained their supremacy in education, business, and the professions, while the remainder of the population retained its sense of "difference" from the dominant peoples and some adopted Roman Catholicism.

Three traditional castes

A further distinction is made between the peoples of the plateau and those of the coast, who are called *côtiers*. The coastal peoples feel deprived of the education, power, and wealth that is concentrated on the plateau. Since independence, the government has been composed of *côtiers*, and a conscious effort has been made to keep the Merina elite of the plateau from power.

CULTURAL LIFE

The cultural milieu. The culture is basically Indonesian. Arabic and Islamic contributions include an intricate system of divination, or *sikidy*, and calendrical features, such as the Arabic-derived names of the days of the week, which also apply to the markets held on those days. The coastal areas of the west, north, and south might be expected to show African cultural elements, but, apart from some Bantu words, these are often difficult to identify conclusively.

The state of the arts. The conquest of the plateau peoples by the French and their subsequent assimilation of Western values have deprived them of most of their traditional institutions. In music, however, Western dance and musical instruments have been adapted to Malagasy rhythms. The tube zither, the conch, and the cone drum are of Indonesian origin, while other types of drums and animal horns suggest African influence. Folk music has been retained, but much of the singing consists of Western church hymns and chants adapted to the distinctive Malagasy musical style.

Social and religious life on the plateau centres upon the church congregation, and the cultural emphasis on ancestral tombs is now largely expressed in Christian terms. More time, money, and care are spent on building tombs than houses. The dead are always brought back to their ancestral tombs, however long or far away they have spent their lives. Tombs are opened every few years, the remains taken out and carried in procession with much ceremony, then replaced after being rewrapped in new shrouds, which are still woven from locally produced silk, coloured with natural, herbal dyes. The male peasants wear cloth trousers with tunics reaching to the knees. Women wear cloth dresses but wrap a silk cloth under one arm and over the other shoulder, even when wearing Western fashions.

The cult of the dead

The coastal peoples have retained more of their traditional customs. Funeral practices are similar to those of the plateau, with local variations of detail. In the eastern forests men wear shorts rather than trousers, and many still wear the short tunic that is woven from raffia fibres. In the far south some older men wear a homespun silk cloth that is wrapped around the waist and between the legs, but most have adopted imported cotton clothes.

The Mahafaly have a remarkable wood-carving industry, and their tombs of coloured stones and carved wooden posts are the most beautiful on the island. The Betsileo also have a thriving wood-carving industry, making in-

Handicrafts laid furniture of valuable hardwoods. They also produce ornamental cloths of very finely woven raffia and have become specialists in the production of coloured straw hats. Betsileo and Merina women have become experts in French-style embroidery, sewing, and dressmaking.

The Malagasy language is rich in proverbs, and there is now an extensive written literature including poetry, legend, history, and scholarly works, as well as contemporary themes. Literary production is aided by an excellent printing industry, for which the Merina have shown a flair since learning it from the London Missionary Society in the 1820s. The peoples of the southeast still preserve their manuscripts in Arabic script with great reverence; few can be more than 200 years old, although some may be copies of much earlier manuscripts.

Cultural institutions. The government encourages the blending of old and new cultural expressions, and a number of new seasonal festivals have been promoted, including the Festival of Rice, the Festival of the Trees, the Festival of the Party, and Independence Day. Towns, churches, schools, and private groups hold concerts or dances, and in the cities there are cultural associations based on the members' home districts.

The main libraries and museums, located in Antananarivo, include the National Library, the Municipal Library, and the National Archive. There are also the library of the Malagasy Academy, the university library, and the university museum. There are museum collections of Malagasy culture and archaeology; natural science collections include a zoo with animals specific to Madagascar.

The press and broadcasting. There are daily and other newspapers published in French and Malagasy, and a government gazette is also published. The island receives radio, television, and telephone service. (A.S.)

For statistical data on the land and people of Madagascar, see the *Britannica World Data* section in the BRITANNICA BOOK OF THE YEAR.

History

EARLY HISTORY

Archaeological investigations of this century indicate that human settlers reached Madagascar about AD 700. Although the huge island lies geographically close to Bantu-speaking Africa, its language, Malagasy, belongs to the distant Malayo-Polynesian language family. There are, nonetheless, a number of Bantu words in the language, as well as some phonetic and grammatical modifiers of Bantu origin. Bantu elements exist in every dialect of Malagasy and appear to have been established for some time.

As a people, the Malagasy represent a unique blend of Asian and African physical and cultural features found nowhere else in the world. Although on the whole Asian features predominate, African ancestry is prevalent and African influences in Malagasy material and nonmaterial culture are widespread. The most plausible theory for this circumstance is that the seafarers of the Malayo-Polynesian world who settled Madagascar initially arrived by way of eastern Africa and the Comoros after these areas had already been colonized by Bantu-speaking Africans. There is also some evidence that Bantu speakers inhabited portions of western Madagascar prior to the 17th century, only to eventually become completely assimilated into Malagasy culture. After the 14th century, important Afro-Arab influences entered Madagascar and spread through much of the island. Apart from the colonial French, settlers from overseas appear to have stopped coming to Madagascar by about 1600.

Madagascar before 1650. Much of Madagascar was populated by internal migration before the beginning of the 16th century, giving the theretofore empty lands their *tompontany* (original inhabitants, or "Masters of the Soil"). Yet politically the island remained fragmented. Most of the nearly 20 ethnic groups that make up the modern Malagasy population did not attain any form of "national" consciousness until new political ideas arrived from abroad in the 1500s and began to spread through the island. A host of written European accounts from the 16th and early 17th centuries fails to reveal any large state or

empire, and few of the Malagasy oral traditions collected since the mid-19th century go back that far in time.

Still, small local states were found at many points along the coast visited by European ships. The capitals were almost always located near river mouths, territorial domains were invariably small, and rulers were independent of one another. Alliances and wars were usually short-lived affairs involving limited economic objectives and little loss of life, and they seldom led to any border adjustments. Economies were pastoral or agricultural, often a mixture of both, and there were no radical differences in wealth. In some areas the rulers appeared to be absolute, while in others elders and priests had the preponderant influence. In one area in southeastern Madagascar, later to become known as Fort-Dauphin (site of the French East India Company fort of that name), early Europeans believed they had found a Muslim state in existence among the Antanosy people of the region. It was ruled by a "Moorish king" and had an aristocracy with privileges deriving presumably from Islām. Their collective name was Zafindraminia, or descendants of Raminia, the ultimate great ancestor.

In the first quarter of the 16th century, Portuguese navigators also reported a number of coastal towns in northern Madagascar that were architecturally similar to Kilwa, a once important entrepôt in what is today Tanzania. The towns belonged to an Afro-Arab commercial network in the western Indian Ocean which undoubtedly predated the 16th century. At the town of Vohemar, once the island's northeastern centre of international trade, the blend of Malagasy and Afro-Arab customs produced an arts and crafts tradition that was quite original.

Portuguese explorers who visited the Matitana River valley in southeastern Madagascar witnessed the arrival of a group of Afro-Arabs ("Moors from Malindi") between 1507 and 1513. Within one or two generations the descendants of this group had intermarried and merged with the local *tompontany* to form another group known as the Antemoro. By the 1630s the Antemoro had formed a theocratic state, which was the only state in Madagascar at the time to possess written texts. Using the Arabic alphabet, the texts were written in the Malagasy language and were both religious and secular in nature. Proximity to Islām became a major criterion among the Antemoro for the right to rule, and there is little doubt that the four Antemoro sacerdotal clans were far closer to the Muslim faith than were the Zafindraminia of the Fort-Dauphin area. In time, Antemoro holy men, traveling far and wide within Madagascar, came to influence other Malagasy in both religion and government.

Political evolution 1650 to 1810. Unknown to the early coastal visitors from Europe, new and historically pivotal dynasties also were beginning to form in southwestern and central Madagascar toward the mid-16th century. Two of them, the Maroserana in the southwest and the Andriana-Merina in central Madagascar, would go on to create vast empires, each with its own apex and decline, between about 1650 and 1896, the year the French annexed Madagascar. While the Maroserana were able to emplace their rulers over several south-central peoples, the most outstanding achievement of the dynasty was the creation of two states in western Madagascar, Menabé and Boina. These states later combined into the Sakalava empire, which controlled most of western Madagascar and several adjacent areas deep inland.

The Sakalava were originally a group of warriors who came into contact with the Maroserana before 1660, the year the Maroserana ruler, King Andriandahifotsy, founded Menabé. Ultimately, "Sakalava citizenship" was extended to hundreds of west-coast clans as the original Sakalava warriors and their descendants intermarried and merged with them. A sense of unity also came from religion, as the Maroserana royals, upon death, became the sacred ancestors of all Sakalava. The Sakalava empire was ultimately weakened by internal power struggles for the throne, by attempts to substitute Islām for the ancestral cult, and, after 1810, by wars with the Merina, a people of the central plateau already on the way to an empire.

The Betsimisaraka confederation, a quasi-state concur-

Small local states

Original settlers

Sakalava empire

rent with the late Sakalava empire, was a brief but successful attempt in the 18th century to unite the coastal peoples of Madagascar's eastern littoral. Ruled by Ratsimilaho, son of an English pirate and a Malagasy princess, the viable confederation extended along more than 200 miles of coastline. After Ratsimilaho's death in 1750, the confederation began an abrupt, though prolonged, disintegration.

Merina
kingdom

The Merina kingdom was founded toward the end of the 16th century in the swampy Ikopa valley on the central plateau. Antananarivo (Tananarive) became its capital. In the 18th century, Imerina was divided among four warring kings. One of them, Andrianampoinimerina, who reigned 1787–1810, reunited the kingdom about 1797. He gave it uniform laws and administration and sold slaves to the French on the coast, using the guns he got in return to conquer his neighbours, the Betsileo. Under Andrianampoinimerina, Merina society was divided into a ruling noble class (Andriana), a class of freemen (Hova), and a slave class (Andevo). At Andrianampoinimerina's death, he left his son a single political ambition: "The sea will be the boundary of my rice field" (*i.e.*, of his kingdom).

Early European contacts. Madagascar is mentioned in the writings of Marco Polo, but the first European known to have visited the island is Diogo Dias, a Portuguese navigator, in 1500. It was called the Isle of St. Lawrence by the Portuguese, who frequently raided Madagascar during the 16th century, attempting to destroy the incipient Muslim settlements there. Other European nations also moved in; in 1642 the French established Fort-Dauphin in the southeast and maintained it until 1674. One of their governors, Étienne de Flacourt, wrote the first substantial description of the island. In the late 17th and early 18th centuries, Madagascar was frequented by European pirates (among them Captain William Kidd) who preyed upon shipping in the Indian Ocean.

In the 18th century, the Mascarene Islands to the east were colonized by the French with the help of Malagasy slaves. Two attempts at fortified settlements failed, one at Fort-Dauphin by the Count de Modave, the other at the Bay of Antongil by Baron Benyowski; however, French trading settlements prospered, notably at Tamatave.

THE KINGDOM OF MADAGASCAR

Formation of the kingdom (1810–61). Andrianampoinimerina's son, Radama I (1810–28), allied himself with the British governor of the nearby island of Mauritius, Sir Robert Farquhar. In order to prevent reoccupation of the east coast by the French, Farquhar supported Radama's annexation of the area by supplying him with weapons and advisers and giving him the title "king of Madagascar." At the same time, Radama agreed to cooperate with Britain's new campaign to end the slave trade. In 1817 he captured the east coast town of Tamatave, from which he launched annual expeditions against the coastal populations. He eventually conquered almost the entire east coast, the northern part of the island, and most of the two large Sakalava kingdoms. Only the south and a part of the west remained independent. The French retained only the small island of Sainte Marie. In addition, Radama invited European workmen, and the London Missionary Society spread Christianity and influenced the adoption of a Latin alphabet for the Malagasy language. Radama died prematurely in 1828; he was succeeded by his widow, Ranaivalona I, who reversed his policy of Europeanization. She expelled Christian missionaries and persecuted Malagasy converts. A few Europeans maintained external trade and local manufacture, but eventually they also were expelled. The British and French launched an expedition against Ranaivalona but were repulsed at Tamatave in 1845. By the time of her death (1861), Madagascar was isolated from European influence.

Merina
expansion

Outside influences (1861–95). Ranaivalona was succeeded by her son, Radama II, who readmitted the foreigners. English Protestants and French Roman Catholics vied for supremacy, while businessmen obtained excessive concessions. This policy led to Radama's overthrow by the Merina oligarchy in 1863. The head of the army, Rainilaiarivony, a Hova, became prime minister and remained in power by marrying three queens in succession: Raso-

herina, Ranaivalona II, and Ranaivalona III. He embarked on a program of modernization, and in 1869 he caused Protestantism to be adopted and suppressed the Malagasy religion. European-style ministries were created and governors set up in the provinces. Villages were supervised by former soldiers. Education was declared obligatory and placed under the direction of the Christian missions. A code of laws was worked out that combined ancient customs with Western practices such as monogamy.

The French began to extend their influence over the Sakalava, and the first "Franco-Merina" war (1883–85) ended with an ambiguous treaty: France was given a settlement at Diégo-Suarez and a resident at Antananarivo, but the institution of a protectorate was temporarily avoided. The succeeding period was marked by disorder and internal strife. In 1890 the British recognized Madagascar as a French protectorate, but Rainilaiarivony refused to submit to French suzerainty. In January 1895, French troops landed at Majunga, and on Sept. 30, 1895, they occupied Antananarivo. The prime minister was exiled. The queen signed a treaty recognizing the protectorate and was maintained on the throne as a figurehead.

The
French
invasion

THE FRENCH PERIOD

The colonial period (1896–1945). French occupation soon extended to the entire part of the island conquered by the Merina. But in Imerina itself, armed guerrilla bands (the Menalamba, or "red togas") resisted modernization and French rule. The French Parliament voted the annexation of the island on Aug. 6, 1896, and sent General Joseph-Simon Gallieni first as military commander, then as governor-general. Slavery was abolished. Gallieni put down the insurrection, subdued the oligarchy, and sent the queen into exile on Feb. 27, 1897. In 1898 the old Merina kingdom was pacified; Gallieni then undertook the difficult task of subjugating the independent peoples. Two insurrections, in the northwest (1898) and in the southeast (1904), were quickly put down and, when he left the island in 1905, unification had been achieved. The Merina governors had been replaced by French administrators with leaders taken from local peoples. The teaching of French in the schools was made compulsory. Customs duties favoured French products, though Malagasy enterprise was also encouraged. The Tamatave-Antananarivo railroad was begun, roads were built, and a modern health service was inaugurated.

The economic development of the island continued under Gallieni's successors. The railroad and its branch lines were completed in 1913. A second line, the Fianarantsoa-Manakara, was finished in 1935. Automobile roads increased after 1920, airlines after 1936. The cities and seaports were built up and equipped, and loans were contracted in France. Exports were confined to agricultural products and raw materials for industry. Rice, cassava (manioc), rubber, raffia, meat, and graphite predominated at first. Between World Wars I and II, coffee, vanilla, cloves, and tobacco, introduced by the Europeans and then taken up by native planters, became more important. Three-quarters of all trade was with France. Material aspects of life became Westernized, especially in the cities, and half the population became Christianized.

In 1915 a nationalist secret society, the Vy Vato Sakelika (VVS), was outlawed. In 1920 a teacher, Jean Ralaimongo, launched a campaign in the press to give the Malagasy "subjects" French citizenship and to make Madagascar a French *département*. When France failed to respond to the demand for assimilation, the movement turned toward nationalism. In 1940 Madagascar, though hesitant at first, rallied to the Vichy government. Then came a blockade, occupation by the British and South Africans (1942), and finally a return to Free France.

Malagasy
nationalism

The French Union (1946–58). In the elections of 1945, two Malagasy nationalists were elected to the French parliament. The constitution of 1946, creating the French Union, made Madagascar an Overseas Territory of the French Republic, with representatives to the Paris assemblies and a local assembly at Antananarivo. Six provincial assemblies were created later. The political struggle erupted into violence on March 30, 1947, with a full-

The
1947–48
rebellion

scale insurrection in eastern Madagascar. The leaders of the Democratic Movement for Malagasy Renewal (Mouvement Démocratique de la Rénovation Malgache), including the three representatives to the French national assembly, were outlawed. While an official count of lives lost in the revolt records about 11,000 dead, it is certain that thousands more of the Malagasy populace perished from famine, cold, and psychological misery while hiding from both the French Army and the insurgents in the island's inhospitable tropical forests.

A period of political inactivity followed until the 1950s. After the "Overseas Territories Law" of 1956 gave Madagascar an executive elected by the local assembly, Vice Premier Philibert Tsiranana founded the Social Democratic Party (Parti Social Démocrate; PSD), which, though most of its members were non-Merina from the coastal areas, offered to cooperate with the Merina. In 1958 France agreed to let its overseas territories decide their own fate. In a referendum on September 28, Madagascar voted for autonomy within the French Community. On Oct. 14, 1958, the autonomous Malagasy Republic was proclaimed; Tsiranana headed the provisional government.

THE MALAGASY REPUBLIC

The opposition regrouped under the name of Congress Party for the Independence of Madagascar (Ankoton'ny Kongresy ny-Fahaleovantena Malagasy), which included both Protestant Merina dissidents and communists. Antananarivo was this party's stronghold; it also had some support in the provinces but, owing to the electoral system established by the PSD, held only three seats in the legislature.

The PSD also settled the provincial question: executive power in the local assemblies was vested in a minister delegated by the central government. Tsiranana was elected president of the republic, and he was instrumental in obtaining its independence on June 26, 1960. Tsiranana and the PSD remained in power until 1972. Under his regime, successive development plans were inspired, according to Tsiranana, by a "grass roots socialism," and were aimed at improving the lot of the peasantry. In foreign policy, the bond with France remained strong, and close relations were established with the United States, West Germany, Taiwan, South Africa, and other anticommunist powers.

Tsiranana was reelected in January 1972, but political and labour unrest, and his own poor health, led him to appoint Major General Gabriel Ramanantsoa as prime minister with full powers of government. A plebiscite on Oct. 8, 1972, confirmed Ramanantsoa as head of government; there was no longer a president after Tsiranana resigned on October 11. The new head of government initiated radically different foreign and domestic policies. Under new agreements with France, French military and naval forces were removed from the island, and all French citizens were to be treated as aliens. Ties were established with the Soviet Union and other communist nations, and the country was withdrawn from the Franc Zone. In 1973, a rural reorganization program—in which elected committees would sell produce to state-owned companies—was initiated, and the government began to take control of joint French-Malagasy organizations.

In the wake of political and social unrest, on Feb. 5, 1975, Ramanantsoa handed power over to a former minister of the interior, Colonel Richard Ratsimandrava. He assumed the titles of president and prime minister but was assassinated six days later. A military directorate was then established; it dissolved on June 15, after naming Lieutenant Commander Didier Ratsiraka president and head of the Revolutionary Council. A referendum on Dec. 21, 1975, approved Ratsiraka as president under a new constitution that set up the Democratic Republic of Madagascar. Ratsiraka, sworn in as president on Jan. 4, 1976, continued the policies begun by Ramanantsoa. He nationalized the banks, insurance companies, and the nation's mineral resources and solidified his nation's ties with the communist powers. The formation in 1976 of the National Front for the Defense of the Revolution—

The
Tsiranana
regime

The
Ratsiraka
regime

a coalition of formerly banned political parties headed by Ratsiraka—further increased the president's control of the government and of Madagascar's political life.

Ratsiraka was reelected without opposition in 1983. Under the banner of scientific socialism he further extended government control over the economy. Despite several loans from the International Monetary Fund to help keep the country's economy afloat, however, it declined drastically. In 1986, Ratsiraka reversed the country's course completely. Laws were altered across-the-board to allow for a free-market economy. In June 1990, France responded by forgiving Madagascar its huge debt of four billion French francs. Overseas investors began to show renewed interest in the country, particularly South Africa. In September 1990 a "new era of cooperation" began between Madagascar and South Africa with a meeting at Antananarivo between Ratsiraka and South African President F.W. de Klerk. Ratsiraka lost the controversial presidential elections held in December 2001. A recount eventually declared Marc Ravalomanana the winner. (H.J.D./R.K.K.)

For later developments in the history of Madagascar, see the BRITANNICA BOOK OF THE YEAR.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, sections 943, 96/11, and 978, and the *Index*.

BIBLIOGRAPHY. General works that provide concise overviews of the geography, history, and economy of the island nation include FREDERICA M. BUNGE (ed.), *Indian Ocean. Five Island Countries*, 2nd ed. (1982); RITA STEVENS, *Madagascar* (1988); and BERNARDINE BAILEY *et al.*, *Madagascar in Pictures*, rev. ed. (1988).

Physical and human geography. *Land and nature:* The natural history of the island is studied in R. BATTISTINI and G. RICHARD-VINDARD (eds.), *Biogeography and Ecology in Madagascar* (1972); M.D. JENKINS (ed.), *Madagascar: An Environmental Profile* (1987); and ALISON JOLLY, PHILIPPE OBERLÉ, and ROLAND ALBIGNAC (eds.), *Madagascar* (1984), the last with a focus on conservation. On animals and plants, see THEODOR HALTENORTH and HELMUT DILLER, *A Field Guide to the Mammals of Africa Including Madagascar* (1980; originally published in German, 1977; also published as *The Collins Field Guide to the Mammals of Africa Including Madagascar*, 1988); IAN TATTERSALL, *The Primates of Madagascar* (1982), a scholarly work with analysis of unique habitats; P. CHARLES-DOMINIQUE *et al.*, *Nocturnal Malagasy Primates: Ecology, Physiology, and Behavior* (1980); OLIVIER LANGRAND, *Guide to the Birds of Madagascar*, trans. from French (1990); and H. PERRIER DE LA BATHIE, *Flora of Madagascar*, trans. from French (1981).

People and culture: CONRAD PHILLIP KOTTAK *et al.* (eds.), *Madagascar: Society and History* (1986); and PIERRE VERIN, *The History of Civilization in North Madagascar* (1986; originally published in French, 1972), are both interdisciplinary studies of the diverse indigenous peoples of the area. Other studies of Malagasy peoples include RICHARD HUNTINGTON, *Gender and Social Structure in Madagascar* (1988); MAURICE BLOCH, *Placing the Dead: Tombs, Ancestral Villages and Kinship Organization in Madagascar* (1971); MARCELLE URBAIN-FAUBLÉE, *L'Art malgache* (1963), on Malagasy art; and LEONARD FOX (ed. and trans.), *Hainteny: The Traditional Poetry of Madagascar* (1990).

Economy: *Madagascar, Recent Economic Developments and Future Prospects* (1980), is a World Bank country study focusing on economic conditions. Broader analysis is found in FREDERIC L. PRYOR, *The Political Economy of Poverty, Equity, and Growth: Malawi and Madagascar* (1990).

Administration and social conditions: For analyses of political situations influencing the development of the area, see MAUREEN COVELL, *Madagascar: Politics, Economics, and Society* (1987); and O. MANNONI, *Prospero and Caliban: The Psychology of Colonization* (1990; originally published in French, 1984), exploring the psychological meaning of colonial dependency of Madagascar for its peoples. (A.I.S.)

History. Developments up to the end of the 18th century are surveyed in RAYMOND K. KENT, *Early Kingdoms in Madagascar, 1500–1700* (1970). HUBERT DESCILAMPS, *Histoire de Madagascar*, 4th rev. ed. (1972), is invaluable as the first historical synthesis. For broader chronological scope leading into the 20th century, see MERVYN BROWN, *Madagascar Rediscovered: A History From Early Times to Independence* (1978); RAYMOND K. KENT, *From Madagascar to the Malagasy Republic* (1962, reprinted 1976); and RAYMOND K. KENT (ed.), *Madagascar in History: Essays from the 1970s* (1979). (R.K.K.)

Madrid

Madrid is the capital of Spain, its largest city, and a national centre of arts and industry. With its surrounding province, also called Madrid, it forms one of the new autonomous regions of the post-Franco era. Its capital status reflects the centralizing policy of the 16th-century Spanish king Philip II and his successors. The choice of Madrid as capital, however, was also the result of Madrid's previous obscurity and neutrality, in that it lacked ties with an established, non-royal power, rather than of any strategic, geographic, or economic considerations. Indeed Madrid is deficient in other characteristics that might be thought to qualify it for a leading role. It is not on a major river, in the way that so many European cities are; the playwright Lope de Vega, referring to a magnificent bridge over the distinctly unimposing waters of the Manzanares, suggested either selling the bridge or buying another river. Madrid does not possess mineral deposits or other natural wealth, nor was it ever a destination of pilgrimages, although its patron saint, San Isidro, enjoys the all but unique distinction of having been married to another saint. Even the city's origins seem inappropriate for a national capital, since its earliest historical role was as the site of a small Moorish fortress on a rocky outcrop—part of the northern defenses of what was then the far more important city of Toledo, 40 miles (65 kilometres) to the southwest.

It was in 1607, a whole generation after Philip II took the court to Madrid in 1561, that Philip III officially made the city the national capital, a status it has retained ever since. Under the patronage of Philip and his successors, Madrid developed into a city of curious contrasts, preserving its old, overcrowded centre, around which developed palaces, convents, churches, and public buildings. This combination has created a quintessentially Spanish city, a true capital in character, with a sparkle and vitality all its own.

This article is divided into the following sections:

Physical and human geography	329
The landscape	329
The city site and climate	
The city layout	
The people	329
The economy	331
Industry, commerce, and finance	
Transportation	
Administration and social conditions	331
Cultural life	331
History	332
The early period	
Development under the Bourbon kings	
Modern Madrid	
Bibliography	333

Physical and human geography

THE LANDSCAPE

The city site and climate. Madrid lies almost exactly at the geographical heart of the Iberian Peninsula. It is situated on an undulating plateau of sand and clay known as the Meseta (derived from the Spanish word *mesa*, or "table") at an altitude of 2,100 feet (635 metres) above sea level, making it one of the highest capitals in Europe. This location, together with the proximity of the Carpetvetónica Range, is partly responsible for the weather pattern of cold, crisp winters accompanied by sharp winds. Sudden variations of temperature are possible, but summers are consistently dry and hot, becoming especially oppressive in July and August, when temperatures sometimes rise above 100° F (37.8° C). Average temperatures range between 41° and 75° F (5° and 24° C), while average

rainfall varies between a low of less than one-half inch (11 millimetres) in July up to about two inches in October, usually the rainiest month of the year. The temperate times of year are spring and fall, which are also the most attractive seasons for visitors.

The city layout. Madrid is a city of contrasting styles, reflecting clearly the different periods in which change and development took place. The old centre, a maze of small streets around a few squares in the vicinity of the imposing Plaza Mayor, contrasts with the stately Neoclassical buildings and grand boulevards created by the most eminent architects of their day. Modern office buildings in the centre and swaths of apartment blocks around the outskirts attest to the styles and economic realities of present-day development. Much of Madrid gives the impression of being cramped. When Madrid was first made the capital, the king obliged the city's inhabitants to let a floor of their houses to ambassadors and visiting dignitaries, which prompted many people to build structures with only one floor or sometimes (in the so-called *casas a la malicia*, or "spite houses") having two floors but with a facade giving the impression of only one. Subsequent development of the city generated an enormous demand for land, particularly with the extensive construction of public buildings and convents. The last of Madrid's four sets of city walls was built in 1625 and was not demolished until 1860 (by which time the population of the city had quadrupled). The situation was not alleviated even when Napoleon's brother Joseph, who briefly interrupted the Bourbon line of kings, demolished the convents to create more open space. Joseph's nickname of "El Rey Plazuelas" ("King of the Small Plazas"; one of the few complimentary ones he was given) derived from the squares he created. They did little to appease the ecclesiastical authorities, whose alienation contributed to his downfall. One of the squares, the Plaza de Oriente, facing the palace of the same name, was cleared of 56 houses, a library, a church, and several convents.

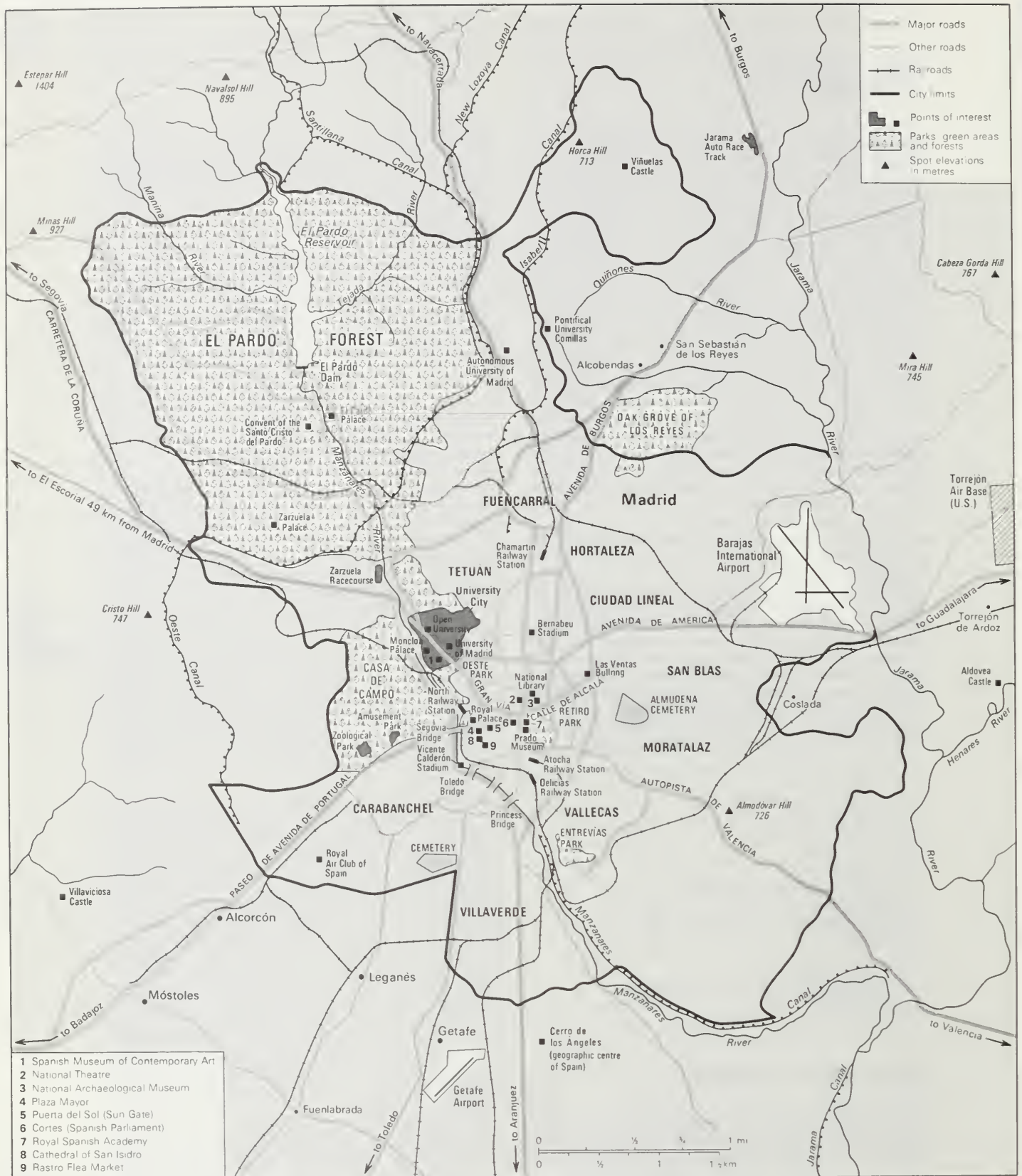
"Los Madriles" ("the Madrids") is a traditional phrase that acknowledges the fact that each *barrio* (quarter) has developed its own style. There was also a geographical and social distinction among the *barrios altos*, *barrios centrales*, and *barrios bajos*. The last, spilling downhill from the Plaza Mayor along the Calle de Toledo toward the river, are still poor, albeit picturesque. Later development, also accommodating Madrid's poorer citizens, spread down toward the reclaimed marshland on both sides of the river, where low-cost housing is still to be found. Just over the brow of the hill is the Rastro, the popular flea market. Despite a number of urban development plans, Madrid did not spread into the open spaces around it, not even crossing the Manzanares River until 1948. By contrast, however, the city as a whole has some extensive parks, with more open space overall than Paris. Some, like El Pardo or Casa de Campo, are survivals of hunting parks; the Retiro, on the other hand, is the site of a former royal palace.

Madrid has not escaped the problems common to so many modern cities. Pollution can be intense, traffic a problem. Personal safety is not as certain as it once was in the days of the *serenos* (night watchmen). But the city has preserved the charm, character, and vivacity that give it and its inhabitants a style of their own—an important aspect of modern Spain, where each region seeks to express its own identity.

THE PEOPLE

The flow of migration to Madrid, attracted chiefly by the city's expanding industrial belt, has created a modern population representative of the entire Spanish nation. A traditional nickname for the Madrileños is *gatos* ("cats"),

Creation of
open space



Madrid.

originally coined in the Middle Ages as a reference to the ability of local troops to scale castle walls. It would be no less apt as a reference to the local life-style and the late hours kept by the city's modern inhabitants, although keeping late hours is also common in other parts of Spain, especially in the heat of summer. People eat late, theatres and cinemas begin performances late as a matter of course, and the siesta is by no means dead, although the introduction of modern business methods and the influx

Late-hour life-style

of foreign interests have tended to lead to the introduction of the *semana inglesa* (literally "English week"), doing away with the long midday break. The city offers a wealth of cultural events and entertainments. It is not a cosmopolitan city as such, and yet its cultivated people tend to be widely read, while the young are up-to-date with the latest pop music. There is no large foreign population, nor are there distinct ethnic groups, although the number of *casas regionales*—regional clubs catering to people who

have come to work from all over the country—reflects the source of labour for Madrid's burgeoning industrial and commercial sectors. Since the early 1970s the Latin-American population has become more numerous as the historic flow of those migrating to the Americas to escape economic misery or political conflict has been reversed. Madrid is a city that, with its style and flair, absorbs and holds those who live there or know it. Its inhabitants have a reputation for being attached to it; in the words of a local proverb, "From Madrid to heaven, and in heaven a little window from which to see it."

THE ECONOMY

Industry, commerce, and finance. Being the centre for government, finance, and insurance is a factor that has for long contributed to the prosperity of the capital, not to mention activities such as tourism and the city's position as the nation's transportation hub. In the postwar period (which, of course, for Spaniards means post-1939—since that was when the Spanish Civil War ended, and since Spain was neutral in World War II) the city became an important manufacturing centre for the automotive and aircraft industries, and for electric and electronic equipment, and optics, the production of plastics and rubber, and consumer goods. Madrid, with Barcelona, dominates publishing in Spain.

Transportation. The road and rail systems both converge on the capital from all corners of the country. A subway system, the Metro, serves Madrid with various lines that extend throughout the city. The international airport of Barajas lies about eight miles east of the city. A motorway (expressway) system encircles Madrid in a roughly pentagonal shape, coming to a point in the south. Other major motorways radiate from the encircling artery in all directions. There are numerous bus routes operated by municipal and private authorities; minibus service is also available. The road-building programs of the 1960s, when much was sacrificed to the convenience of automobile owners, have since been recognized as less than wholly beneficial. Some of the overpasses introduced to speed up traffic flow have since been dismantled.

ADMINISTRATION AND SOCIAL CONDITIONS

With the return of democracy to Spain in the late 1970s and the development of autonomous regional governments, more emphasis has been placed on local consultation and issues such as the future of the environment. In

1982 the city administration carried out a massive public opinion survey to find out what people really wanted at the neighbourhood level. The resulting General Ordinance Plan (Plan General de Ordenación) was an attempt to establish a long-term, full-scale scheme for future directed growth, aiming not only to modernize the infrastructure of essential services but also to improve the quality of life in the city. Despite the introduction of the autonomous regions, Madrid continues to be the focus of Spain's government. It is also a bishopric, headquarters of the army corps, and residence of the captain general of the first military region. The supreme court and government ministries are also located there, besides the Cortes (the Spanish parliament), which is housed in a 19th-century Neoclassical building. Characterized by the bronze lions (made from melted down Moorish cannon) flanking the entrance, it is one of the smallest parliament buildings in Europe. Local administration is under the direction of a mayor and city council.

CULTURAL LIFE

Modern pressures have perhaps inhibited the extensive street life for which Madrid has been famous, although people still live very much in the streets, especially during the intense heat of summer when the café terraces fill and people stroll up and down in the evenings. Modern culture, in the form of film, theatre, and music, is extensively represented, as is to be expected in a city with several major universities and academies. But the cultural life for which Madrid was once noted, the *tertulias*, that is to say the informal conversational gatherings and informal societies, have all but faded, along with the elegant cafés that housed them. Madrid's literary traditions, its associations with Lope de Vega, Calderón, Cervantes, Quevedo, Pérez Galdós, Larra, Baroja, and Azorín, continue in the city's varied cultural life, as demonstrated by the fact that it is one of the major publishing centres for the Spanish-speaking world.

Modern Madrid has attractions at all levels. The 23,000-seat bullring, Las Ventas, is the largest in Spain and the place where the novice bullfighters have to display their skills in the *alternativa* (the occasion on which a matador kills his first bull) in order to become established. The season runs from March to October. There are two major football (soccer) teams, with the annual matches against the Barcelona teams among the high points of the year. Important matches are played in two stadia, the

The General Ordinance Plan

Bull-fighting tradition

By courtesy of the Spanish National Tourist Office



Palacio Real and gardens, Madrid.

Postwar industrial development

125,000-seat Santiago Bernabéu and the 70,000-seat Vicente Calderón. The *verbenas*, special fiestas held in each quarter in honour of its patron saint, are regular public events, especially in warm weather, with San Isidro (mid-May) taking pride of place. The *zarzuelas* (light opera of mildly satirical flavour, indulging in topical comment and set by tradition in Madrid) are commonly held in the open air at this time. There are, in fact, more than 40 parks and public gardens, the principal ones being the Retiro, Campo del Moro, Casa de Campo, and Oeste Park, not to mention the curious temple of Debod (an ancient Egyptian temple acquired by Spain at the time of the construction of the Aswān High Dam) near Rosales, with splendid panoramic views over the city.

Museums abound. Unusual ones include those for theatre, railways, and (understandably enough for Spain) *taurromaquia*, the bullfight. The city is richly endowed with artistic masterpieces: tapestries in the Casa de Cisneros (the mayor's residence) and the Royal Palace (Palacio Real); paintings by Bruegel and Titian in the convent of the Descalzas Reales; and Spanish and foreign masters in the Palacio de Liria, home of the dukes of Alba. The most famous collection is housed in the Prado, which displays the artworks collected by the Spanish monarchy over the ages and reflects the pattern of Spain's alliances. Charles V and Philip II were patrons of Venetian art; Philip IV was a great collector in the 17th century; and the accession of the Bourbon family led to an influx of French works. Spain's control of the Netherlands led to a solid Flemish section. El Casón del Buen Retiro is nearby and houses 19th- and 20th-century works, the most famous probably being Picasso's "Guernica." This painting in 1981 was sent to Spain from New York City in accordance with Picasso's directive that the painting be moved there only after democracy had returned to the country.

Notable among an abundance of libraries are the prestigious National Library (Biblioteca Nacional) and the Library of the Royal Palace, acclaimed for its historic collection. Madrid is also famous for its secondhand bookshops, and the Feria del Libro (book fair), held in the spring, is a widely heralded event.

Madrid is Spain's foremost centre of higher education and includes several of the country's leading universities, among which are the Open University (Universidad Nacional de Educación a Distancia), the Complutensian University, and the Polytechnical University, all in University City, and to the north the Autonomous University.

History

The early period. The Arab town, or medina, grew around the alcazar (castle) on a promontory overlooking the Manzanares River. The name Majerit first appears in AD 932, when the Christian king Ramiro II of León razed the town's walls, but there are traces of earlier (even prehistoric) habitation. The Christian king Alfonso VI of Castile and León captured the town from the Muslims in 1083, and thereafter a number of kings of Castile spent time there. The parliament (Cortes) was called there as early as 1309. The alcazar was damaged in an earthquake in 1466 and the subsequent medieval palace was extended by various monarchs, notably Charles I and Philip II. In this period the town grew to the east up both sides of what are now the Calle Mayor and the Calle de Segovia, with the Moors (who continued to live there until after the Christian reconquest of Spain was completed in 1492) jammed into the southwest corner, which is still called the Moreria. The whole of the city at this time was only 500 by 900 yards in area. Some of the street patterns of the pre-16th century city remain, but few buildings; one that still stands is the much-restored Casa de los Lujanes, where it is believed the French king Francis I was once held prisoner. Charles I enjoyed hunting near Madrid, and it is said that the widening of the city gates to accommodate his carriages opened the cramped streets to heavier traffic, a process that was to increase when the court was properly established in the city. By 1598 the population of Madrid had reached 60,000, and by the time of the first extant plan of Madrid (Pedro Teixeira's in 1656) it

had grown to be an imposing city of 100,000 people and 11,000 buildings.

Under the Habsburg monarchs, the Madrid de los Austrias expanded even more rapidly. The foreign ministry (1634), the Casa de Cisneros, and the Segovia Bridge date from this time, as does the church of San Isidro el Real. (It is also, incidentally, still Madrid's temporary cathedral—Nuestra Señora de la Almudena, begun rather tardily in 1883, is still under construction.) Architects such as Juan de Herrera (who died in 1597) and Francisco de Mora contributed to the monumental quality of the city. But the most striking contribution of this period is generally considered to be the Plaza Mayor, designed by Juan Gómez de Mora and built between 1617 and 1619; it was modified after the great fire of 1790. Graceful in concept, it is surrounded by five-story houses with balconies and topped with steeples. Nine archways open onto the plaza at oblique angles from surrounding streets, and the continuous arcade at street level contains shops and restaurants. Bullfights (in those days conducted by noblemen on horseback), fireworks displays, and plays all took place there, as did the grim ceremonies of the Inquisition. Until 1765 public garrotings were also carried out there. The last bullfight to take place in the Plaza Mayor was in honour of the wedding of Isabella II in 1846.

Development under the Bourbon kings. No less impressive a landmark dates from the next great phase of the city's growth, under the Bourbons, whose side Madrid took against the Habsburgs in the War of the Spanish Succession (1701–14), although the city was briefly occupied by pro-Habsburg troops. The Royal Palace was begun by Philip V after the disastrous fire that destroyed the Alcazar on Christmas night, 1734. His grandiose plan, with 23 inner courts, was never realized, although the finished work did have 500 rooms. It was a fitting addition to the other major city features created under his patronage, the Royal Spanish Academy (Real Academia Española), the National Library, and the Royal Academy of History (Real Academia de Historia). The Royal Palace, with its elegant granite and limestone walls, contains a ceiling by Tiepolo in the throne room and, in the Armeria, one of the world's finest collections of armour, including the swords of the conquistadores Hernán Cortés and Francisco Pizarro. The last king of Spain actually to live there was Alfonso XIII, whose apartments have been preserved just as he left them when he abdicated in 1931. The royal family now resides in the more private and less ostentatious La Zarzuela Palace, set in its own grounds to the northwest of Madrid. A less important palace was also favoured by the Bourbon kings—the Palacio del Buen Retiro; its gardens, which were much admired for their French style, are still a much frequented open space.

The greatest Bourbon builder was Charles III (1759–88), who is known as the mayor-king for his interest in the growth and development of the city and his many contributions to its skyline. This was the age of Enlightenment, and Charles tried to establish a harmony between leisure and science, culture and industry. His style was cosmopolitan, reflecting the tastes of the Europe of his time. With his concern for the appearance of the city, its gates, avenues, and trees, he anticipated the designs of modern city planners. He relied heavily in his schemes on the works of three Neoclassical architects, Francisco Sabatini, Ventura Rodríguez, and Juan de Villanueva. During this period the city continued to grow eastward to the present Plaza de la Independencia, which is the site of a monumental arch, the Puerta de Alcalá, built in 1778 and still a key landmark. Another famous landmark, the central post office, in the Puerta del Sol, also dates from this time. All distances in the country are still measured from the zero-kilometre stone beneath its wall. The square itself holds a firm place in people's affections; it is particularly popular on New Year's Eve when people go to hear the clock strike midnight and to eat the 12 grapes that supposedly ensure good luck in the following months. In earlier times the square was significant because of the stagecoaches that left from there to all parts of the peninsula. Its appeal made it the site of various innovations in urban amenities, from the first gas lamps in 1830 through the first mule

Habsburg rule

Growth under the Bourbon kings

Charles III, the mayor-king

The Prado

Christian takeover

trams and first public urinal to the first electric streetlights and electric streetcars. Charles III would doubtless have approved of all these innovations. In his time he set up the Botanic Garden (which still exists), with a "physic" garden from which anybody could (and still can) collect medicinal herbs. He planned a natural history and science museum next to it but died before it could be completed.

Madrid was occupied by French troops during the Napoleonic Wars, and Napoleon's brother Joseph was installed on the throne. On May 2, 1808, there was a mass uprising against Joseph, leading to what the Spaniards term "la Guerra de la Independencia." Ferdinand VII, on his return in 1814 from imprisonment by Napoleon, bestowed the title of "heroic" upon the city. In 1819 the building intended by Charles III to house a natural history and science museum was completed. Into it Ferdinand moved artworks of the royal collection, until then scattered among various palaces. This was the start of what was to become one of the world's major art galleries, the Prado. The Madrid of this period can still be studied in close detail, thanks to the remarkable model constructed by León Gil Palacios in 1830. It was during this period that the city expanded to the north, under the direction of Joaquín Vizaíno, a nobleman who was also mayor (as was customary at the time). He is also known as the man who introduced such innovations as street numbers for buildings, street lighting, and municipal refuse collection. The Paseo del Prado was extended by two new boulevards, the Paseo de los Recoletos and the Paseo de la Castellana. For years this area presented an almost rural atmosphere and featured imposing town houses with great gardens. Now it is the site of tall office blocks, apartment buildings, luxury hotels, and embassies, as well as the National Library, the National Archaeological Museum, the Queen Sofia Art Centre, and the Thyssen-Bornemisza Museum. One of the area's mansions (now a bank) belonged to the Marqués de Salamanca, who in 1872 also contributed to the drive northward by building 28 streets on a grid plan, starting from the Calle de Alcalá and parallel to the Paseo. Still bearing his name, this area remains one of the most elegant *barrios* in Madrid.

Modernization under the Plan Castro

Somewhat earlier, in 1860, the Plan Castro, also referred to as the *Ensanche*, or widening, had further expanded and modernized the city, adding convenience and meeting the economic and commercial needs of the time. It was the first comprehensive, forward-looking modern plan for Madrid. However, it was to be frustrated by population growth, land speculation, and the poor areas that sprang up outside the planned zones.

Modern Madrid. During the 20th century, Madrid grew dramatically and became a modern commercial and administrative centre. This transformation had begun by 1910, when a major new landmark appeared. The *barrio* of San Bernardo was bisected by a broad way running from the Calle de Alcalá downhill to the Plaza de España, where the city's first high-rise commercial buildings were erected. This, the Gran Vía, was designed to be the main street of the city, and it has a characteristic vitality, with cinemas, coffeehouses, shops, and banks.

With the advent of the Republic in 1931, radial and ring roads appeared, and the Paseo de la Castellana was extended even farther. The city suffered heavily in the Civil War, with two years of aerial and artillery bombardment and combat lines drawn up as close as University City. War damage to public buildings was repaired, and ambitious reconstruction plans were drawn up, but these were by and large unsuccessful in practice. The city spread out-

ward, swallowing its own suburbs; between 1948 and 1951 Madrid's city limits were extended to cover a total of 205 square miles (531 square kilometres), about 10 times its previous area. A period of land speculation and uncontrolled urban sprawl followed. During the 1950s and '60s, the city expanded rapidly into the surrounding countryside. Industrial districts developed south and east of the historic centre, encircled by areas of ramshackle, low-rise squatter housing and poorly built high-rise public housing for the workers who streamed into the city. Meanwhile, elite suburban developments of single-family homes proliferated on Madrid's northern and western peripheries. High-rise office buildings spread north from the historic centre along the Paseo de la Castellana and displaced many of the elegant mansions that had lined the street.

During the 1980s and '90s, Madrid enjoyed renewed prosperity and rising living standards. New suburban developments of single-family and two-family houses appeared on the city's outskirts, even in the traditionally working-class south and east. At the same time, an extensive renovation program improved the quality of older working-class neighbourhoods and restored many aging buildings in the historic centre. Legislation passed toward the end of the 20th century protected many of the city's most prized structures and promoted the preservation of less-famous historic buildings. At the beginning of the 21st century, Madrid's planners worked to balance present-day prosperity with respect for the city's rich historic legacy.

BIBLIOGRAPHY

General works: *All Madrid*, trans. from Spanish (1979); and ALASTAIR BOYD, *The Companion Guide to Madrid and Central Spain* (1974), are descriptive guides. JUAN ANTONIO CABEZAS, *Madrid*, 3rd ed. (1971), examines all districts and cultural institutions of the capital. ARCHIBALD LYALL, *Well Met in Madrid* (1960); and ELENA SAINZ, *Living in Madrid* (1981), are more personal accounts.

History: FEDERICO BRAVO MORATA, *Historia de Madrid*, 4 vol. (1966-78); and JOSÉ AMADOR DE LOS RÍOS, *Historia de la villa y corte de Madrid*, 4 vol. (1860-64; reprinted 1978), give wide coverage of historical events. Modern history, and especially the Civil War, is presented in GEORGE HILLS, *Battle for Madrid* (1976); MATILDE VÁZQUEZ and JAVIER VALERO, *La guerra civil en Madrid* (1978); DAVID JATO MIRANDA, *Madrid, capital republicana* (1976); DAN KURZMAN, *Miracle of November: Madrid's Epic Stand, 1936* (1980); and JOSÉ MANUEL MARTÍNEZ BANDE, *Frente de Madrid* (1976).

Social and economic development: NINA EPTON, *Madrid* (1964), describes the inhabitants of the various quarters of the city; MICHAEL KENNY, *A Spanish Tapestry: Town and Country in Castile* (1969), compares life in Madrid to that in a rural parish; DAVID R. RINGROSE, *Madrid and the Spanish Economy, 1560-1850* (1983); JULIO VINUESA ANGULO, *El desarrollo metropolitano de Madrid* (1976); and SANTOS JULIA DÍAZ, *Madrid, 1931-1934: de la fiesta popular a la lucha de clases* (1984), are scholarly studies of social developments and conflicts.

Art and monuments: HARRY B. WEHLE, *Great Paintings from the Prado Museum* (1963), includes a sketch of the museum's history; F.H. SANCHEZ CANTON, *The Prado*, trans. from the French, new rev. ed. (1966), also includes a historical essay; MANUEL LORENTE, *The Prado, Madrid*, 2 vol. (1965), is an impressive collection of reproductions; CONSUELO LUCA DE TENA and MANUELA MENA, *Guide to the Prado* (1980; originally published in Spanish, 1980), is a more recent survey. JONATHAN BROWN and J.H. ELLIOTT, *A Palace for a King: The Buen Retiro and the Court of Philip IV* (1980), is a well-researched description of the pleasure palace and its notable art collection. See also AUREA DE LA MORENA BARTOLOMÉ *et al.*, *Catálogo monumental de Madrid* (1976). (B.E./T.J.Co.)

Malta

An independent republic, Malta (Repubblika Ta' Malta) comprises a small but strategically important group of islands in the central Mediterranean Sea. Throughout a long and turbulent history, the archipelago has played a vital role in the struggles of a succession of powers for domination of the Mediterranean and in the interplay between emerging Europe and the older cultures of Africa and the Middle East. As a result, Maltese society was molded by centuries of foreign rule, with influences ranging from Arab to Norman to English.

There are five islands—Malta (the largest), Gozo, Comino, and uninhabited Kemmalett (Comminotto) and Filfla—lying some 58 miles (93 kilometres) south of Sicily, 180 miles (290 kilometres) north of Libya, and about 180 miles east of Tunisia, at the eastern end of that constricted portion of the Mediterranean Sea separating Italy from the African coast. The islands cover a combined land area of 122 square miles (316 square kilometres). Valletta is the capital, although Birkirkara is the largest city.

The article is divided into the following sections:

Physical and human geography 334

- The land 334
 - Relief
 - Drainage and soils
 - Climate
 - Plant and animal life
 - Settlement patterns
- The people 335
- The economy 336
 - Industry and tourism

- Agriculture and fishing
- Finance
- Transportation
- Government and social conditions 336
 - Government
 - Education
 - Health and welfare
- Cultural life 337
- History 337
- Bibliography 338

Physical and human geography

THE LAND

Relief. Malta Island measures about 17 miles at its longest distance from southeast to northwest and about 9 miles at its widest distance from east to west. The main physical characteristic of Malta is a well-defined escarpment that bisects it along the Victoria Lines Fault running along the whole breadth of the island from Point ir-Raġeb (west of Nadur Tower) to the coast northeast of Gharghur. The highest areas are coralline limestone uplands that constitute a triangular plateau. Ta' Zuta (829 feet [253 metres]), to the west. The uplands are separated from the surrounding areas by blue clay slopes, while undercliff areas are found where the coralline plateau has fallen and forms a subordinate surface between the sea and the original shore. The total shoreline is 85 miles.

To the north the escarpment is occasionally abrupt and broken by deep embayments. To the south, however, the plateaus gradually descend from about 600–800 feet into undulating areas of globigerina (derived from marine protozoa) limestone less than 400 feet high. On the west are deeply incised valleys and undercliff areas, while on the east are several valleys that descend to the central plains.

The west coast of Malta presents a high, bold, and generally harbourless face. On the east, however, a tongue of high ground known as Mount Scerberras separates the bays of Marsamxett and Grand Harbour. These deepwater

harbours contribute to the strategic importance of Malta. They are associated with nine seasonal creeks that include those of Sliema, Lazzaretto, Msida, and Newport. The northern shore is again bare and craggy, characterized by its coves and hills, which are separated by fertile lowlands.

In Gozo the landscape is characterized by a broken coralline plateau to the north and by low-lying globigerina limestone plains and hills to the south. The highest point, in the west, is 578 feet. The total shoreline is 27 miles.

Drainage and soils. Malta possesses favourable conditions for the percolation and underground storage of water. The impermeable blue clays provide two distinct water tables between the limestone formations. The principal source for the public supply of water has been the main sea-level water table. The absence of permanent streams or lakes and a considerable loss of rainfall, however, have made water supply a problem. This problem has been combated with an intensive reverse-osmosis desalination program. About 70 percent of Malta's daily water needs are supplied by desalination plants throughout the islands.

Maltese soils are mainly young or immature and thin. By law, when soils are removed from construction sites, they must be taken to agricultural areas, and level stretches in quarries are often covered with carted soil. Organic refuse from the towns is also used. Consequently, the soils are unusual and are partly a manufactured medium.

Climate. The climate is typically Mediterranean, with hot, dry summers, warm and sporadically wet autumns,

Shortage
of fresh
water

Limestone
and coral



The Dockyard complex in Grand Harbour, Malta, showing (from foreground) French Creek, Senglea, Dockyard Creek, Bormla (Cospicua) at the head of Dockyard Creek, Fort St. Angelo on Vittoriosa point, Kalkara Creek, Kalkara, and Rinella Creek.

Financial Times, London—Robert Harding Picture Library

MAP INDEX

Cities and towns

Attard	35 53 N 14 27 E
Birkirkara (Birchircara)	35 54 N 14 28 E
Birżebbuga (Birzebbugia)	35 50 N 14 32 E
Bormla (Cospicua)	35 53 N 14 31 E
Citta Vecchia, see Mdina	
Cospicua, see Bormla	
Dingli	35 52 N 14 23 E
Għajnsielem (Għain Silem)	36 02 N 14 17 E
Għarb	36 04 N 14 13 E
Għargħur	35 55 N 14 27 E
Għaxaq	35 51 N 14 31 E
Gżira	35 54 N 14 30 E
Fhamrun	35 53 N 14 29 E
Kerċem	36 02 N 14 14 E
Kirkop	35 51 N 14 29 E
Lija	35 54 N 14 27 E
Luqa	35 52 N 14 29 E
Marsa	35 53 N 14 30 E
Marsaskala	35 52 N 14 34 E
Marsaxlokk	35 50 N 14 33 E
Mdina (Citta Vecchia or Medina Notabile)	35 53 N 14 24 E
Mellieħa	35 58 N 14 22 E
Mgarr	35 55 N 14 22 E
Mosta	35 55 N 14 26 E
Mqabba	35 50 N 14 28 E
Msida	35 54 N 14 29 E
Nadur	36 02 N 14 18 E
Naxxar	35 55 N 14 27 E
Notabile, see Mdina	
Paola, see Raħal Ġdid	
Paula, see Raħal Ġdid	
Pawla, see Raħal Ġdid	
Qala	36 02 N 14 19 E
Qormi	35 53 N 14 28 E
Qrendi	35 50 N 14 27 E
Rabat	35 53 N 14 24 E

Rabat (Victoria)	36 03 N 14 14 E
Raħal Ġdid (Paola, Paula, or Pawla)	35 52 N 14 30 E
St Julian, see San Ġiljan	
St Paul's Bay, see San Pawl il-Baħar	
San Ġiljan (St Julian)	35 55 N 14 29 E
San Pawl il-Baħar (St. Paul's Bay)	35 57 N 14 24 E
Sannat	36 02 N 14 15 E
Siggiewi	35 51 N 14 26 E
Sliema	35 55 N 14 30 E
Tarxien (Tarshin)	35 52 N 14 31 E
Valletta (Valetta)	35 54 N 14 31 E
Victoria, see Rabat	
Xewkija	36 02 N 14 16 E
Żabbar	35 52 N 14 32 E
Żebbuġ	35 52 N 14 26 E
Żebbuġ	36 04 N 14 14 E
Żejtun	35 51 N 14 32 E
Żurrieq	35 50 N 14 28 E

Physical features and points of interest

Aħrax, Point Ta' I-	36 00 N 14 22 E
Bajda Ridge	35 57 N 14 22 E
Bengħisa, Point Ta'	35 49 N 14 33 E
Buskett, park	35 52 N 14 24 E
Comino (Kemmuna), island	36 01 N 14 20 E
Comminotto, see Kemmunett	
Dalam Cave, historical site	35 50 N 14 32 E
Dawwara, Point Id-	35 52 N 14 21 E
Filfla, island	35 47 N 14 24 E
Gaulus, see Gozo	
Ġgantija, historical site	36 03 N 14 16 E
Għawdex, see Gozo	



Gozo (Gaulus or Għawdex), island	36 03 N 14 15 E
Grand Harbour	35 54 N 14 31 E
Faġar Qim, historical site	35 50 N 14 27 E
Kemmuna, see Comino	
Kemmunett (Comminotto), island	36 00 N 14 19 E
Malta, island	35 55 N 14 25 E
Marfa Ridge	35 59 N 14 21 E
Marsamxett Harbour	35 54 N 14 30 E
Mediterranean Sea	35 55 N 14 15 E
Mellieħa Bay	35 59 N 14 22 E
Nadur Tower, mountain	35 54 N 14 22 E
Qammieħ, Point Il-	35 58 N 14 19 E
Raħeb, Point Ir-	35 54 N 14 20 E
St Paul's Bay	35 57 N 14 24 E
Wahx, Point Il-	35 57 N 14 20 E
Wardija Ridge	35 56 N 14 23 E

and short, cool winters with adequate rainfall. Nearly three-fourths of the total annual rainfall of about 20 inches (508 millimetres) falls between October and March; June, July, and August are normally quite dry.

The temperature is very stable, the annual mean being 64° F (18° C) and the monthly averages ranging from 54° F (12° C) to 88° F (31° C). Winds are strong and frequent; the most common are the cool northwesterly (the *majjistral*), the dry northeasterly (the *grigal*, or gregale), and the hot humid southeasterly (the *xlokk*, or sirocco). The relative humidity is consistently high and rarely falls below 40 percent.

Plant and animal life. While wild vegetation is sparse, there is an abundance of cultivated potatoes, *sulla* (a leguminous fodder crop), onions, tomatoes, and vines. The forest cover, though poor, features a wide variety of trees, including the carob, fig, and chaste. The Maltese government has initiated a major tree-planting program to improve forestation on the islands. The developed seashore vegetation includes golden and rock samphires, sea campion, spurge, saltwort, and marram grass.

Typical of the few native mammals are the hedgehog, the least weasel, the water and white-toothed shrews, and the pipistrelle and other bats. Rats, mice, and some rabbits also are found. Resident birds include the spectacled and Sardinian warblers, the Manx and Cory's shearwaters, and the blue rock thrush. Linnets, tree and Spanish sparrows, buntings, rock doves, and several species of owl also form breeding populations, while birds of passage include ospreys, rollers, swallows, cuckoos, bee-eaters, and vultures. Common insects are beetles, grasshoppers, flies, mosquitoes, moths, bees, wasps, cockroaches, ants, and several species of butterflies. Ladybirds migrate from Sicily.

Settlement patterns. Until the mid-19th century, the Maltese lived mainly in the relative seclusion of clustered

villages and hamlets; the fragmentation of farmholdings accentuated the individuality of the farming community. The *zuntier*, or church square, was the traditional focus of village life. With the growth of the Dockyard complex in the latter part of the 1800s, new settlements appeared around Grand Harbour, and the Sliema metropolitan region developed in the 20th century into the most fashionable part of Malta. The advent of industrial estates near major villages somewhat stemmed the exodus from the rural areas. Higher living standards have given rise to residential developments all over the island; its central and northwest areas are now densely populated. Overbuilding has been a cause for serious concern, spawning legislation meant to contain the ecological threats thus posed.

Gozo conserves its own rural character. The architecture of the development at Ta' Ċenċ successfully blends with the island's natural beauty and is aesthetically stimulating. Comino is more rural still, with only a handful of residents, no cars, and two hotels.

THE PEOPLE

The islands' ethnic and linguistic composition reflects the heritage of many rulers. A European atmosphere predominates as a result of close association particularly with southern Europe. About 95 percent of the islanders are Maltese-born, and the remainder includes mostly persons of English and Italian descent. During the 20th century, the increasing rate of Anglo-Maltese marriages added a new dimension to the ethnic structure of the population.

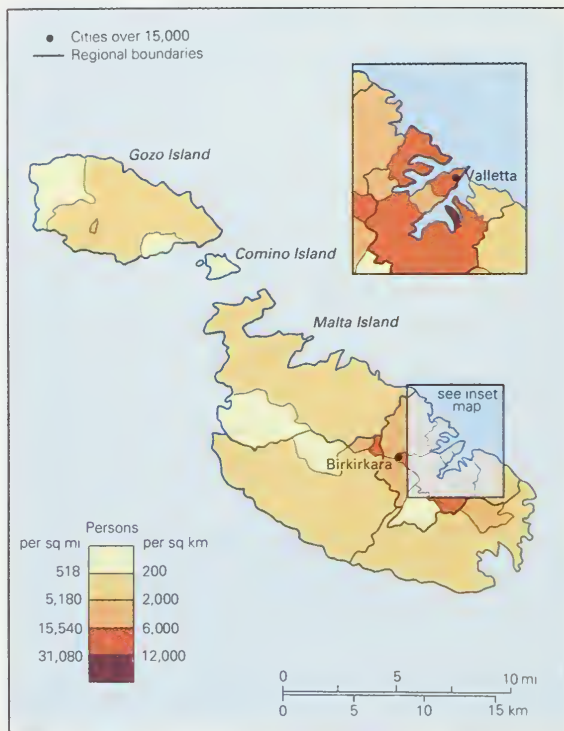
Maltese—the medium of daily conversation—is a distinct language that resulted from the interaction and fusion of North African Arabic and a Sicilian form of Italian. It is the only Semitic language that is officially written in Latin script; it became an official language of Malta in 1934. English, the other official language, is the medium

Concern with overbuilding

of instruction in the schools. Italian was the language of church and government until 1934 and is still understood by a sizable sector of the population.

Although Roman Catholicism is the state religion, there is full freedom for all religious beliefs. The islands are an independent province of the church, with two dioceses at Malta and Gozo and two bishops serving the cathedrals at Valletta and Rabat (Victoria). There are two Roman Catholic cathedrals, at Mdina and Valletta, and an Anglican cathedral, at Valletta. There is a mosque at Corradino Heights.

Malta has one of the highest population densities in the world. The population has somewhat stabilized, however, with a considerable decline in the birthrate since the 1950s. The death rate has remained fairly stable, having fallen only slightly, while the infant mortality rate has dropped significantly. Emigration, formerly encouraged and even financed by the government, has also tapered off; the primary destination of emigration is Australia. The majority of the country's immigrants are repatriates. The age structure of the population is fairly evenly distributed, and the life expectancy is about 74 years.



Population density of Malta.

THE ECONOMY

Malta's only exploited mineral resource is the globigerina limestone that is used as building stone, although the country has offshore reserves of petroleum. Its other assets are its deep harbours, an adaptable, skilled labour force, and a strategic position as both a fueling centre and, until March 1979, a military and naval base. The economy, therefore, has been somewhat artificial and until 1979 determined by the vicissitudes of war and peace in the Mediterranean. In the 1950s Britain began the retrenchment of its naval and military forces on Malta, a move that necessitated a drastic diversification of the economy. A series of five- and seven-year plans were supported by government grants, loans, and other fiscal incentives to encourage private investment.

Industry and tourism. Economic plans professed to build on a tripod basis of industry, agriculture, and tourism. In fact, however, industrial growth lagged behind these plans, resulting in the successful establishment of only a few multinational corporations (mainly producing textiles). From 1971 the government increasingly took over weak enterprises, sometimes closing them. Since 1987 new development has concentrated on manufacture of industrial components, including computer parts, in-

struments, and other high-tech goods, as well as a large variety of consumer products (toys, cosmetics, detergents, processed foods) and more traditional goods such as lace, silver filigree, pottery, glassware, and canework. Foreign investment in manufacturing is encouraged and facilitated by the Malta Development Corporation.

Tourism is a major source of income. The influx of tourists and some immigration spurred the building of hotels and housing, but there are questions about the islands' capacity to cater to an annual total of tourists greater than the country's population. Besides shipbuilding and transshipment services, the establishment of backup facilities for oil companies and of other outlets for traditional Maltese skills in cross-cultural dealings produced jobs and foreign currency earnings that boosted the Maltese currency, making it one of the strongest in the world.

Agriculture and fishing. Agricultural development is hampered by infertile soils and the lack of adequate water supplies, and large amounts of food and beverages are imported. Most farming is carried out on small terraced strips of land that preclude the introduction of large-scale mechanization. Consequently, agriculture continues to decline in terms of cultivated land. The farming labour force has grown increasingly older, female, and part-time, but agricultural production has risen gradually because of improved techniques in the cultivation of some crops, especially horticultural ones. The major crops are grains, vegetables, fruit (especially citrus), and fodder for domestic consumption and early potatoes and onions for export. Flowers, seeds, plants, and cuttings are exported as well. Cattle, pigs, sheep, poultry, and goats are also raised.

Fishing is seasonal and not fully developed. The fishing boats are small, and the catch is affected by weather conditions. The fish population off the island is quite sparse and in the late 1980s began to be supplemented by agriculture. During the summer months, small lamps are used by fishermen (lampara fishing) to attract the fish, mainly the very popular lampuka, or dolphin (*Coryphaena hippurus*).

Finance. The Central Bank and the Malta Development Corporation were both founded in 1968. Also in that year, Malta joined the International Monetary Fund. The Malta Export Trade Corporation was founded in 1989. In 1971 the island entered into a special trade relationship with the European Communities (EC), a relationship that was subsequently revised with the aim of rendering it more favourable to Malta in its delicate stage of economic development. Malta is trying to secure full membership status in the EC.

Transportation. The road system connects all towns and villages and includes a coast road, a panoramic road, and a regional road, which connects Msida to the coast road via St. Andrews. Bus services radiating from Valletta provide inexpensive and frequent internal transportation. Taxis are abundant, and most families have a private automobile. There is no railway. Several daily sailings connect Malta and Gozo, and, in addition to a regular ferry and car-ferry service, Malta and Sicily are connected by fast commercial catamarans. Scheduled airlines, including the national airline (Air Malta, based at Luqa Airport), connect Malta with most European capitals as well as with Africa, especially North Africa, and the Middle East.

GOVERNMENT AND SOCIAL CONDITIONS

Government. The 1964 Independence Constitution, under which Malta was a constitutional monarchy and parliamentary state, was amended in 1974 to make Malta a republic within the Commonwealth. Its head of state is a president appointed by the Maltese Parliament, which is elected by universal adult suffrage for a term of five years and is basically derived from the British model. Local features include a single chamber with 69 members, while election is by proportional representation from 13 electoral divisions. An amendment adopted in 1987 guarantees a majority of seats to a party receiving more than 50 percent of the total votes in the general election. The two major parties are the Nationalist Party and the Malta Labour Party. The president acts on the advice of the Cabinet, which consists of the prime minister and other ministers (some assisted by parliamentary secretaries) and

is collectively responsible to Parliament. There is no municipal government in the islands.

Maltese law, which was codified mainly during the period from 1854 to 1873, is largely based on the Napoleonic Code and Napoleonic law. Procedural common law and some commercial and maritime affairs are regulated by English principles, but judiciary precedent is not binding. Maltese is the language of the courts. Civil and criminal jurisdiction is almost exclusively vested in the Superior Court and the Court of Magistrates. The chief justice and judges of the Superior Court are appointed by the president on the advice of the prime minister, and their duties are apportioned throughout the court system. The magistrates, who are appointed in the same way, sit in the lower courts.

Between 1964 and 1972, Malta's main defense dispositions were those contained in a 1964 Anglo-Maltese defense agreement, with the United Kingdom guaranteeing mutual assistance. Since then, through a constitutional amendment, Malta has followed a policy of neutrality and nonalignment and maintains its own regular armed forces.

Education. The Labour government radically altered the education system, which was previously structured on British models and strongly influenced by the Roman Catholic church. Compulsory education was extended to include all children from the ages of 6 to 16. An attempt at establishing an extreme form of the "comprehensive" system was abandoned; streaming (the grouping of students by age and intellectual ability) and examinations were at first discarded but later reintroduced; purely technical institutes were not compelled to follow the program. At the tertiary level, a student-worker scheme was introduced in 1978, students working for six months and studying for six months, thereby linking admission to institutions of higher learning to the availability of employment. This system was largely revoked by the Education Act of 1987, and admission to institutions of higher learning is now based completely on competence.

The University of Malta, founded as a Jesuit college in 1592 and established as a state institution in 1769, was refounded in 1988. It offers courses in most disciplines and has a prestigious medical school. Its modern campus at Tal-Qroqq also houses the International Maritime Law Institute and the Mediterranean Academy of Diplomatic Studies. The historic Old University building in Valletta is now the seat of the university-linked Foundation for International Studies and its associated bodies, the International Environment Institute, the Mediterranean Institute, and the Euro-Mediterranean Centre on Marine Contamination Hazards (created by the Council of Europe). Malta is also the site of the Regional Marine Pollution Emergency Response Centre for the Mediterranean Sea, operated jointly by the International Maritime Organization (IMO) and the United Nations Environment Programme (UNEP).

Health and welfare. A scheme to integrate health services, essentially consisting of free hospital care and some domestic assistance, was introduced after 1971. Certain measures deemed to be necessary preliminaries to a wider-ranging national health scheme gave rise to conflict with the local medical association, and heavy reliance was then placed on foreign doctors. Since 1988 Malta has been the seat of the United Nations International Institute on Aging, and thus special attention is given to geriatrics.

In 1956 social insurance was introduced to cover employees more than 14 years of age and self-employed or unemployed persons between 19 years and pension age. A comprehensive contributory insurance scheme was introduced in 1971, integrating a variety of earlier legislation. The plan included a pension amounting to two-thirds of an individual's salary at the time of retirement. The first step toward comprehensive national health insurance was taken in 1979 with the introduction of free hospitalization.

CULTURAL LIFE

Malta's cultural influences stem largely from its history of foreign domination and the predominance of the Roman Catholic church. Folk traditions have evolved mainly around the festa to celebrate the patron saint of a village, marked by processions and fireworks. Good Friday also

is celebrated with colourful processions in several villages. *Imnarja*, the Feast of St. Peter and St. Paul, which takes place on June 29, is the principal folk festival; it is highlighted by folksinging (*ghana*) contests and fried-rabbit picnics at Buskett. The annual Carnival is celebrated at Valletta with vigorous dancing displays that include the Parata, a sword dance commemorating the Maltese victory over the Turks in 1565, and Il-Maltija, the Maltese national dance. Soccer is the most popular sport in Malta, and Ta' Qali National Stadium is the site of important local and international matches.

Apart from its unique Neolithic ruins near Raħal Ġdid (Paola, Paula, or Pawla) and Tarxien, Malta contains important examples of its flourishing architectural school of the 17th and 18th centuries. It was essentially Classical with a balanced overlay of Baroque decorations. The Italian artists Caravaggio and Mattia Preti spent several years in Malta, the latter's most important paintings embellishing many of Malta's churches.

In the 20th century, a vernacular architecture was developed by Richard England and others. The composer Charles Camilleri introduced folk themes into his works, while Maltese literature was enriched by the poetry of the national bard, Dun Karm. An interesting theatrical upsurge led by John Schranz paralleled the emergence of Francis Ebejer as a brilliant playwright. Alfred Chircop and Luciano Micallef have gained prominence with their abstract paintings, while Gabriel Caruana has excelled in ceramics.

Valletta is the centre of many of Malta's cultural institutions: the National Museum of Archaeology, the National Museum of Fine Arts, the War Museum, the Manoel Theatre (one of Europe's oldest theatres still in operation), and the Foundation for International Studies. The National Library of Malta dates from the late 18th century and houses a large collection as well as the archives of the Knights Hospitalers. The Folk Museum and the Museum of Political History are located at Vittoriosa. Until the early 1990s, Maltese radio and television stations had been operated exclusively by the Malta Broadcasting Authority, but a change in legislation has opened the way for privately operated broadcasting stations. There are two daily newspapers in Maltese and one in English.

For statistical data on the land and people of Malta, see the *Britannica World Data* section in the BRITANNICA BOOK OF THE YEAR.

History

The earliest archaeological remains date from about 3800 BC. Neolithic farmers lived in caves like those at Dalam (near Birżebbuġa) or villages like Skorba (near Nadur Tower) and produced pottery that seems related to that of contemporary eastern Sicily. An elaborate cult of the dead of Stone Age or Copper Age culture evolved about 2400 BC. Initially centring around rock-cut collective tombs such as those at Ġgantija (near Xagħra) and Haġar Qim (near Żurrieq), it culminated—probably through contacts with the cultures of the Cyclades and Mycenae—in the unique underground burial chamber (hypogeum) at Hal Saflieni (near Raħal Ġdid). This culture came to a sudden end about 2000 BC, possibly as a result of invasions. The culture that replaced it, of southern Italian flavour, is evidenced today only by fragmentary remains. Bronze Age tools and weapons have been found at Borġ in-Nadur (near Birżebbuġa) and Tarxien Cemetery (near Raħal Ġdid), while Iron Age relics from about 1200 to 800 BC include cart ruts at Bingemma (near Nadur Tower).

Between the 8th and 6th centuries BC, contact was made with Semitic cultures. Evidence is scanty, however, and a few inscriptions found on Malta constitute the only indication of a Phoenician presence. There is more substantial proof of the Carthaginian presence in the 6th century BC; coins, inscriptions, and several Punic-type rock tombs have been found. It is certain that in 218 BC Malta came under Roman political control. Originally part of the praetorship of Sicily, the islands were subsequently given the status of municipium and were thus allowed to coin their own money, send ambassadors to Rome, and

The
judicial
system

The
arts

Folk
traditions
and
festivals

Domina-
tion by
Carthage
and Rome

control domestic affairs. According to tradition, St. Paul was shipwrecked in Malta in AD 60 and began to convert the inhabitants. The Maltese have been Christians uninterruptedly since that time.

With the division of the Roman Empire, in AD 395, Malta was given to the eastern portion dominated by Constantinople. Until the 15th century, it followed the more immediate fortunes of nearby Sicily, being ruled successively by Arabs (who left a strong effect on the language), Normans (who advanced the legal and governmental structures), and a succession of feudal lords. In 1530, however, it was ceded to the Order of the Hospital of St. John of Jerusalem (the Knights Hospitalers), a religious and military order of the Roman Catholic church. Malta became a fortress and, under the Knights' grand master, Jean de La Valette, successfully withstood the Ottoman siege of 1565. The new capital city of Valletta became a town of splendid palaces and unparalleled fortifications. Growing in power and wealth—owing mainly to their maritime adventures against the Turks—the Knights left the island an architectural and artistic legacy. Although there was little economic and social contact between them and the Maltese, they managed to imprint their cosmopolitan character on Malta and its inhabitants.

In 1798 Napoleon Bonaparte captured the island, but the French presence was short-lived, and the Treaty of Amiens returned the island to the Knights in 1802. The Maltese protested and acknowledged Great Britain's sovereignty, subject to certain conditions incorporated in a Declaration of Rights. The constitutional change was ratified by the Treaty of Paris in 1814.

British rule

Malta's political status under Britain underwent a series of vicissitudes in which constitutions were successively granted, suspended, and revoked. British demands for Malta's military facilities dominated the economy, and the Dockyard became the colony's economic mainstay.

The island flourished during the Crimean War and was favourably affected by the opening of the Suez Canal in 1869. Self-government was granted in 1921 on a dyarchical basis whereby Britain shared power and responsibility with Maltese ministers who were elected by the legislature. But the battle over the relative roles of the English, Italian, and Maltese languages took its toll, and in 1936 the islands reverted to a strictly colonial regime in which full power rested in the hands of the governor. During World War II, the islands repelled the Axis powers against severe odds, having been one of the most heavily bombed targets of that conflict. As a result, it was awarded the George Cross, Britain's highest civilian decoration. Self-government was granted in 1947, revoked in 1959, and then restored in 1962. Malta finally achieved independence within the Commonwealth on Sept. 21, 1964. It became a republic on Dec. 13, 1974.

The withdrawal of British military and naval personnel from its famous Dockyard—associated with the achievement of independence from the United Kingdom in 1964—created economic and political problems for Malta. From 1964 to 1971, Malta was governed by the Nationalist Party, whose attitude was firm alignment with the West. In 1971, however, when the Malta Labour Party came to power, its policy was nonalignment and special friendship with China and Libya. In 1979 the total closure of the British base and the end of the British alliance were celebrated by the Maltese government as the arrival of "real" independence; however, the existence of serious problems, especially with regard to the guarantee of full and productive employment and the deepening division between local political parties, was acknowledged. The Nationalists were returned to power in 1987 with a policy of seeking full membership in the European Community and transforming the economy into a modern, technologically oriented structure. Malta's growing international importance was acknowledged when it was chosen as the site of the first summit meeting between Soviet President Mikhail Gorbachev and U.S. President George Bush in 1989.

(S.Bu.)

For later developments in the history of Malta, see the BRITANNICA BOOK OF THE YEAR.

BIBLIOGRAPHY

Geography. *General works:* WALTER KÜMMERLY *et al.*, *Malta: Isles of the Middle Sea* (1965); and HARRY LUKE, *Malta: An Account and an Appreciation*, 2nd ed., rev. and enlarged (1960), are illustrated descriptive works with maps. Focus on local landscape is found in ROBIN BRYANS, *Malta and Gozo* (1966); HARRISON LEWIS, *A Guide to the Remote Paths and Lanes of Ancient Malta* (1974); and DOUGLAS LOCKHART and SUE ASHTON, *Landscapes of Malta, Gozo, and Comino* (1989). BODO NEHRING, *Die Maltesischen Inseln* (1966), provides a more scholarly geographic survey. Travelers' guides, some repeatedly revised in well-known publishers' series, include BRYAN BALLS and RICHARD COX, *Traveller's Guide to Malta: A Concise Guide to the Mediterranean Islands of Malta, Gozo, and Comino*, 4th ed. (1981); INGE SEVERIN, *See Malta & Gozo: A Complete Guide with Maps and Gazetteer*, rev. ed. (1984); and PETER MCGREGOR EADIE, *Malta and Gozo*, 3rd ed. (1990).

People: *Minor Islands of the Mediterranean, Gozo, Malta* (1981) is a UNESCO survey of the settlement patterns in the region. JEREMY BOISSEVAIN, *Hal-Farrug: A Village in Malta* (1969, reissued as *A Village in Malta*, 1980), offers a study of social life and customs. The role of religion is explored in JEREMY BOISSEVAIN, *Saint and Fireworks: Religion and Politics in Rural Malta* (1965); and MARIO VASSALLO, *From Lordship to Stewardship: Religion and Social Change in Malta* (1979).

Government and social conditions: PAUL CASSAR, *Medical History of Malta* (1964), is a detailed survey of the development of the essential social service. BARRY YORK, *Malta: A Non-Aligned Democracy in Mediterranean* (1987), offers a short overview of modern politics; and, for an equally brief look at the civil rights, see *Human Rights in Malta* (1985), a report of the International Helsinki Federation for Human Rights.

Economy: Surveys of the economic conditions are presented in M.M. METWALLY, *Structure and Performance of the Maltese Economy* (1977); and R. COHEN, M. MINOGUE, and J. CRAIG, *Small Island Economies* (1983). SALVINO BUSUTTIL, *Devaluation in Malta* (1968), examines the currency question in the period of independence. JOHN C. GRECH, *Threads of Dependence* (1978), explores the problems of dependency on foreign technology. LINO BRIGUGLIO, *The Maltese Economy: A Macroeconomic Analysis* (1988); and MICHAEL FRENDO and JOSEF BONNICI, *Malta in the European Community: Some Economic & Commercial Perspectives* (1989), are later analyses.

Cultural life and institutions: On the architecture, see LEONARD MAHONEY, *A History of Maltese Architecture: From Ancient Times up to 1800* (1988); J. QUENTIN HUGHES, *The Building of Malta During the Period of the Knights of St. John of Jerusalem, 1530-1795*, rev. ed. (1967, reissued 1986); and CHARLES KNEVITT, *Connections: The Architecture of Richard England, 1964-84* (1984). Other arts are discussed in RICHARD ENGLAND (ed.), *Contemporary Art in Malta* (1973), including essays on music, painting, and drama; DANIEL MASSA (ed.), *Individual and Community in Commonwealth Literature* (1979); OLIVER FRIGGIERI, *Storia della letteratura maltese* (1986), an Italian translation from the Maltese of an analysis of Maltese poetry; and, on painting, MARIO BUHAGIAR, *The Iconography of the Maltese Islands, 1400-1900* (1988).

History. General historical surveys are presented in ERIC GERADA-AZZOPARDI, *Malta: An Island Republic* (1979); BRIAN BLOUET, *The Story of Malta*, 3rd rev. ed. (1981); and MARIO BUHAGIAR (ed.), *Proceedings of History Week 1983* (1984). For the earliest periods, see J.D. EVANS, *The Prehistoric Antiquities of the Maltese Islands* (1971). The Middle Ages are studied in ANTHONY T. LUTTRELL (ed.), *Medieval Malta: Studies on Malta Before the Knights* (1975); ERNLE BRADFORD, *The Great Siege: Malta 1565* (1961, reissued 1979); ROSE G. KINGSLEY, *The Order of St. John of Jerusalem: Past and Present* (1918, reprinted 1978); with the history continued in RODERICK CAVALLIERO, *The Last of the Crusaders: The Knights of St. John and Malta in the Eighteenth Century* (1960). The modern period is explored in HENRY FRENDO, *Party Politics in a Fortress Colony: The Maltese Experience* (1979); R. DE GIORGIO, *A City by an Order* (1985); ERNLE BRADFORD, *Siege: Malta to 1940-1943* (1985); GEORGE HOGAN, *Malta: The Triumphant Years, 1940-43* (1978); CHARLES A. JELLISON, *Besieged: The World War II Ordeal of Malta, 1940-1942* (1984); DENNIS AUSTIN, *Malta and the End of the Empire* (1971); J.J. CREMONA, *An Outline of the Constitutional Development of Malta Under British Rule* (1963); EDITH DOBIE, *Malta's Road to Independence* (1967); and HENRY FRENDO, *Malta's Quest for Independence: Reflections on the Course of Maltese History* (1989).

(S.Bu.)

Mammals

The class Mammalia is a group of vertebrate animals, collectively known as mammals, in which the young are nourished with milk from special secreting glands (mammary glands) of the mother. In addition to these characteristic milk glands, mammals are distinguished by several other unique features. Hair is a typical mammalian feature, although in many whales it has secondarily disappeared except in the fetal stage. The mammalian lower jaw is hinged directly to the skull, instead of through a separate bone (the quadrate) as in all other vertebrates. A chain of three tiny bones transmits sound waves across the middle ear. A muscular diaphragm separates the heart and the lungs from the abdominal cavity. Only the left aortic arch of the primitive fourth pair persists (in birds the right arch persists; in reptiles, amphibians, and fishes both arches are retained). Mature red blood cells in all mammals lack a nucleus; all other vertebrates have nucleated red blood cells.

Except for the monotremes (echidnas and duck-billed

platypuses), which lay eggs, all mammals are viviparous (*i.e.*, bear live young). In the placental mammals (including man) the young are carried within the mother's womb, reaching a relatively advanced stage of development before being born. In the marsupials (kangaroos, opossums, and allies) the newborn, incompletely developed at birth, continue to develop outside the womb, attaching themselves to the female's body in the area of her mammary glands. Some marsupials have a pouchlike structure or fold, the marsupium, that shelters the suckling young.

The class Mammalia is worldwide in distribution. It has been said that mammals have a wider distribution and are more adaptable than any other single class of animals, with the exception of certain less complex forms such as the arachnids and insects. This versatility in exploiting the Earth is attributed in large part to the ability of mammals to regulate their body temperatures and internal environment both in excessive heat and aridity and in severe cold.

(Ed.)

The class Mammalia 339

General features 339

Importance to man

Natural history

Reproduction

Behaviour

Ecology

Locomotion

Food habits

Form and function

The skin and hair

Dentition

Skeleton

Muscles

Viscera

Evolution and classification 347

The evolution of the mammalian condition

Classification

Distinguishing taxonomic features

Annotated classification

Critical appraisal

Major mammal orders 353

Monotremata (platypus, echidnas [spiny anteaters]) 353

Marsupialia (kangaroos, bandicoots, phalangers, opossums, koala, wombats) 356

Insectivora (shrews, moles, hedgehogs, tenrecs, solenodons) 364

Chiroptera (bats) 370

Primates (lemurs, lorises, tarsiers, monkeys, apes, and hominids, including man) 377

Edentata (armadillos, sloths, anteaters) 392

Lagomorpha (rabbits, hares, pikas) 398

Rodentia (rats, mice, beavers, squirrels, guinea pigs, capybaras) 401

Carnivora (cats, dogs, bears, skunks, seals, walrus) 412

Cetacea (whales, dolphins, porpoises) 429

Proboscidea (elephants) 434

Sirenia (dugong, manatees) 438

Perissodactyla (horses, asses, zebras, tapirs, rhinoceroses) 439

Artiodactyla (pigs, goats, sheep, cattle, giraffes, deer, camels) 447

Bibliography 457

THE CLASS MAMMALIA

General features

The evolution of the Mammalia has produced tremendous diversity in form and habits. Living kinds range in size from tiny shrews, weighing but a few grams, to the largest of all animals that has ever lived, the blue or sulphur-bottom whale, which reaches a length of more than 30 metres (100 feet) and a weight of 136,000 kilograms (150 tons). Every major habitat has been invaded by mammals that swim, fly, run, burrow, glide, or climb.

There are approximately 4,000 species of living mammals, arranged in about 120 families and 20 orders. The rodents (order Rodentia) are the most numerous of existing mammals, both in number of species and number of individuals, and are one of the most diverse of living lineages. In contrast, the order Tubulidentata is represented by a single living species, the aardvark. The Proboscidea (elephants) and Perissodactyla (horses, rhinoceroses, and allies) are examples of orders in which far greater diversity occurred in mid- and late-Tertiary times (from about 3,000,000 to about 30,000,000 years ago) than today.

The tropical continental areas of the world around are the places of greatest present-day diversity of mammals, although members of the order occur on (or in seas ad-

acent to) all major land masses and on many oceanic islands (the latter principally, but by no means exclusively, reached by bats). Major regional faunas can be identified; these resulted in large part from evolution in comparative isolation of stocks of early mammals that reached these areas. South America (Neotropics), for example, was separated from North America (Nearctic) from Paleocene through much of Pliocene times (about 2,500,000 to 65,000,000 years ago) by inundation of the Panamanian Portal and adjacent areas in Middle America. Evolution of mammalian groups that had reached South America before the break between the continents, or some that "island-hopped" in after the break, proceeded quite independently from that of relatives that remained in North America. Some of the latter became extinct as the result of competition with more advanced groups, whereas those in South America flourished, some radiating to the extent that they have successfully competed with invaders since the rejoining of the two continents. Australia provides a parallel case of early isolation and adaptive radiation of mammals (monotremes, marsupials) thus isolated, but differs in that it was not later connected to any other land mass. The more advanced mammals of the infra-class Eutheria that reached Australia (rodents, bats) did

Radiation
in South
America

so, evidently by "island-hopping," long after the adaptive radiation of the early isolates.

In contrast, North America and Eurasia (Palearctic) are separate land masses but have closely related faunas, the result of connections several times in the Pleistocene and earlier across the Bering Strait. Their faunas frequently are thought of as representing not two distinct units, but one, related to a degree that a single name, Holarctic, is applied to it.

IMPORTANCE TO MAN

Wild and domesticated mammals are so interlocked with man's political and social history that it is impractical to attempt to assess the relationship in precise economic terms. Throughout his cultural evolution, for example, man has been dependent on other mammals for a significant portion of his food and clothing. Domestication of mammals helped to provide a source of protein for ever-increasing human populations and provided means of transportation and heavy work as well. Today, domesticated strains of the house mouse, European rabbit, guinea pig, hamster, gerbil and other species provide much-needed laboratory subjects for the study of human-related physiology, psychology, and a variety of diseases from dental caries to cancer. Recent emphasis on the study of nonhuman primates (monkeys and apes) has opened broad, new areas of research relevant to man's welfare. The care of domestic and captive mammals is, of course, the basis for the practice of veterinary medicine.

Some primitive peoples still depend on wild mammals as a major source of food, and many different kinds (from fruit bats and armadillos to whales) regularly are captured and eaten. On the other hand, the hunting, primarily for sport, of various rodents, lagomorphs, carnivores, and ungulates supports a multibillion-dollar enterprise. In the United States alone, for example, it is estimated that more than 2,000,000 deer are harvested annually by licensed hunters.

Geopolitically, the quest for marine mammals was responsible for the charting of a number of areas in both Arctic and Antarctic regions. The presence of terrestrial fur bearers, particularly beaver and several species of mustelid carnivores (*e.g.*, marten and fisher), was one of the principal motivations for the opening of the American West, Alaska, and the Siberian taiga. Sale of wild-taken furs still is an important industry in both the Old and New Worlds, but ranch-raised animals such as the mink, fox, and chinchilla now have become an important part of the fur industry, which directly and indirectly accounts for many millions of dollars in revenue each year in North America alone.

Aside from pelts and meat, special parts of some mammals regularly have been sought for their special attributes. The horns of rhinoceroses are used in concocting potions in the Orient; ivory from elephants and walrus is highly prized; and ambergris, a substance regurgitated by sperm whales, has been widely used as a base for perfumes (although now mostly replaced by synthetic substitutes).

Some mammals are directly detrimental to man and his activities. Murid rodents (house rats and mice of Old World origin) now occur in virtually all of the world's urban areas and each year cause substantial damage and economic loss. In some rural areas, herbivorous mammals eat or trample planted crops and compete with livestock for food, and native carnivores prey on domestic herds. These and other situations have led man to spend large sums annually to control populations of "undesirable" wild mammals, a practice long deplored by conservationists. Mammals are important reservoirs or agents of transmission of a variety of diseases that afflict man, such as plague, tularemia, yellow fever, rabies, leptospirosis, hemorrhagic fever, and Rocky Mountain spotted fever. The annual "economic debt" resulting from mammal-borne diseases of man and his domestic stock is incalculable.

Many large mammals that competed directly with man for food or space, or were specially sought by him for some reason, have been extirpated entirely or exist today only in parks and zoos; others are in danger of extinction, and their plight is receiving increased attention by a number of

conservation agencies. Perhaps at least some can be saved. One of the most noteworthy cases of direct extirpation by man is the Steller's sea cow (*Hydrodamalis gigas*). These large (up to 12 metres, or 40 feet, long), inoffensive, marine mammals evidently lived in Recent times only along the coasts and shallow bays of the Commander Islands in the Bering Sea. Discovered in 1741, they were easily killed by Russian sealers and traders for food, their meat being highly prized, and the last known live individual was taken in 1768.

Of final note is the esthetic value of wild mammals and the relatively recent predilection of man to expend both considerable energy and resources to study and, if possible, conserve vanishing species, to set aside natural areas where native floral and faunal elements can exist in an otherwise highly agriculturalized or industrialized society, and to establish modern zoological parks and gardens. Such outdoor "laboratories" attract millions of visitors annually and will provide means by which present and future generations of humans can appreciate and study, in small measure at least, other kinds of mammals.

NATURAL HISTORY

The hallmarks of the mammalian level of organization are advanced reproduction and parental care, behavioral flexibility, and endothermy (the physiological maintenance of a relatively constant body temperature independent of that of the environment, allowing a high level of activity). Within the class, ecological diversity has resulted from adaptive specialization in food-getting, habitat preferences, and locomotion.

Throughout the past 70,000,000 years, the mammals have been the dominant animals in terrestrial ecosystems and important in nonterrestrial communities as well. The earliest mammals were small, active, predacious, and terrestrial or semiarboreal. From this primitive stock the mammals have radiated into a wide spectrum of adaptive modes against the background of the diverse environment of the Cenozoic (the last 65,000,000 years). Branches of the ancestral terrestrial stock early exploited the protection and productivity of the trees, whereas other lineages added further dimensions to the mammalian spectrum by adapting to life beneath the ground, in the air, and in marine and freshwater habitats.

Reproduction. In reproductively mature female mammals an interaction of hormones from the pituitary gland and the ovaries produces a phenomenon known as the estrous cycle. Estrus, or "heat," typically coincides with ovulation, and during this time the female is receptive to the male. Estrus is preceded by proestrus, during which ovarian follicles mature under the influence of a follicle-stimulating hormone from the anterior pituitary. The follicular cells produce estrogen, a hormone that stimulates proliferation of the uterine lining, or endometrium. Following ovulation, in late estrus, the ruptured ovarian follicle forms a temporary endocrine gland known as the corpus luteum. Another hormone, progesterone, secreted by the corpus luteum, causes the endometrium to become quiescent and ready for implantation of the developing egg (blastocyst), should fertilization occur. In members of the infraclass Eutheria (known as placental mammals) the placenta, as well as transmitting nourishment to the embryo, has an endocrine function, producing hormones that maintain the endometrium throughout gestation.

If fertilization and implantation do not occur, a phase termed metestrus ensues, in which the reproductive tract assumes its normal condition. Metestrus may be followed by anestrus, a nonreproductive period characterized by quiescence or involution of the reproductive tract. On the other hand, anestrus may be followed by a brief quiescent period (diestrus) and another preparatory proestrus phase. Mammals that breed only once a year are termed monestrous, and exhibit a long anestrus; those that breed more than once a year are termed polyestrous. In many polyestrous species the estrous cycle ceases during gestation and lactation (milk production), but some rodents have a postpartum estrus, mating immediately after giving birth.

The menstrual cycle of higher primates is derived from the estrous cycle, but differs from it in that when pro-

Extinction
of Steller's
sea cow

Influence
on
exploration

gestosterone secretion from the corpus luteum ceases, in the absence of fertilization, the uterine lining is sloughed. In anthropoids other than man a distinct period of "heat" occurs around the time of ovulation.

Egg-laying by monotremes

Monotremes lay shelled eggs, but the ovarian cycle is similar to that of other mammals. The eggs are predominantly yolk (telolecithal), like those of reptiles and birds. Young monotremes are altricial (*i.e.*, in a relatively early stage of development and dependent upon the parent), reaching sexual maturity in about one year. (For additional information see below in *Monotremata*, subhead *Life cycle and reproduction*.)

The reproduction of marsupials differs from that of placentals in that the uterine wall is not specialized for the implantation of the embryos. The period of intrauterine development varies from about eight to 40 days. After this period the young migrate from the vagina to attach to the teats for further development. The pouch, or marsupium, is variously structured. Many species, such as kangaroos and opossums, have a single, well-developed pouch; in some phalangruids the pouch is compartmented, with a single teat in each compartment. The South American caenolestids or rat opossums have no marsupium. The young of most marsupials are dependent on maternal care for considerable periods, 13 to 14 weeks in *Didelphis*. Young koalas are carried in the pouch for nearly eight months, kangaroos to 10 months. (See below in *Marsupialia*, subheads *Life cycle* and *Reproductive adaptations*.)

Reproductive patterns in placental mammals are diverse, but in all cases a secretory phase is present in the uterine cycle and the endometrium is maintained by secretions of progesterone from the corpus luteum. The blastocyst implants in the uterine wall. Villi are embedded in the lining of the uterus. The resulting complex of embryonic and maternal tissues is a true placenta. The uterine lining may be shed with the fetal membranes as "afterbirth" (a condition called deciduate) or may be resorbed by the female (nondeciduate). Placentae have been classified on the basis of the relationship between maternal and embryonic tissues. In the simplest nondeciduate placental arrangement the chorionic villi are in contact with uterine epithelium (the inner surface layer). In the "intimate deciduous" types, seen in primates, bats, insectivores, and rodents, the capillary endothelium (the layer containing minute blood vessels) of the uterine wall breaks down, and chorionic epithelium is in direct contact with maternal blood. In advanced stages of pregnancy in rabbits even the chorionic epithelium is eroded and the embryonic endothelium contacts the maternal blood supply. In no case, however, is there actual exchange of blood between mother and fetus; nutrients and gases must still pass through the walls of the fetal blood vessels.

The period of intrauterine development, or gestation, varies widely among eutherians, generally depending on the size of the animal, but influenced by the number of young per litter and the condition of young at birth. The gestation period of the domestic hamster (*Mesocricetus auratus*) is about two weeks, whereas that of the blue whale is 11 months and of the African elephant 21 to 22 months.

At birth the young may be well developed and able to move about at once (precocial) or they may be blind, hairless, and essentially helpless (altricial). In general, precocial young are born after a relatively long gestation period and in a small litter. Hares and many large grazing mammals bear precocial offspring. Rabbits, carnivores, and most rodents bear altricial young.

After birth young mammals are nourished by milk secreted by the mammary glands of the female. The development of milk-producing tissue in the female mammae is triggered by conception, and the stimulation of suckling the newborn prompts copious lactation. In therians (marsupials and placentals) the glands open through specialized nipples. The newborn young of marsupials are unable to suckle and milk is "pumped" to the young by the mother.

Milk consists of fat, protein (especially casein), and lactose (milk sugar), as well as vitamins and salts. The actual composition of milk of mammals varies widely among species. The milk of whales and seals is some 12 times

as rich in fats and four times as rich in protein as that of domestic cows but contains almost no sugar. Milk has meant an efficient energy source for the rapid growth of young mammals; the weight at birth of some marine mammals doubles in five days.

Behaviour. *Social behaviour.* The dependence of the young mammal on its mother for nourishment has made possible a period of training. Such training permits the nongenetic transfer of information between generations. The ability of young mammals to learn from the experience of their elders has allowed a behavioral plasticity unknown in any other group of organisms and has been a primary reason for the success of the mammals. The possibility of training is one of the factors that has made an increase in neural complexity selectively adaptive. Increased associational potential and memory extends the possibility of learning from experience, and the individual can make adaptive behavioral responses to environmental change. Individual response to short-term change is far more efficient than genetic response.

Some kinds of mammals are solitary, except for brief periods when the female is in estrus. In other species, however, social groups are present. Such groups may be reproductive, defensive, or may serve both functions. In those cases that have been studied in detail a more-or-less strict hierarchy of dominance prevails. Within the social group, maintenance of the hierarchy may depend on physical combat between individuals, but in many cases stereotyped patterns of behaviour evolve to displace actual combat, conserving energy, while maintaining the social structure.

Sexual dimorphism (a pronounced difference between sexes) frequently is extreme in social mammals. In large part this is because dominant males tend to be those that are largest or best armed. Dominant males also tend to have priority in mating, or may even have exclusive responsibility for mating within a "harem." Rapid evolution of secondary sexual characteristics can take place in a species with such a social structure even though the reproductive potential may be low.

A complex behaviour termed "play" frequently occurs between siblings, between members of an age class, or between parent and offspring. Play extends the period of maternal training and is especially important in social species, providing an opportunity to learn behaviour appropriate to the maintenance of dominance.

Territoriality. That area covered by an individual in his general activity is frequently termed the home range. A territory is a part of the home range defended against other members of the same species. As a generalization it may be said that territoriality is more important in the behaviour of birds than of mammals, but data for the latter are available primarily for diurnal species. The phenomenon of territoriality may be more widespread than is supposed. Frequently territories of mammals are "marked," either with urine or with secretions of specialized glands, as in lemurs, a form of territorial labelling less evident to humans than the singing or visual displays of birds. Many mammals that do not maintain territories per se nevertheless will not permit unlimited crowding and will fight to maintain individual distance. Such mechanisms result in more economical spacing of individuals over the available habitat.

Ecology. *Response to environmental cycles.* Mammals may react to environmental extremes by acclimatization, compensatory behaviour, or physiological specialization. Physiological responses to adverse conditions include torpidity, hibernation (in winter), and estivation (in summer). Torpidity may occur in the daily cycle or during unfavourable weather; short-term torpidity generally is economical only for small mammals that can cool and warm rapidly. The body temperature of most temperate-zone bats drops near that of the ambient air whenever the animal sleeps. The winter dormancy of bears at high latitudes is an analogous phenomenon and cannot be considered true hibernation.

True hibernation involves physiological regulation to minimize the expenditure of energy. The body temperature is lowered and breathing may be slowed to as low

Hibernation

Composition of milk

as 1 percent of the rate in an active individual. There is a corresponding slowing of circulation and typically a reduction in the peripheral blood supply. When the body temperature nears the freezing point, spontaneous arousal occurs, although other kinds of stimuli generally elicit only a very slow response. In mammals that exhibit winter dormancy (such as bears, skunks, and raccoons), arousal may be quite rapid. Hibernation has evidently originated independently in a number of mammalian lines, and the comparative physiology of this complex phenomenon is only now beginning to be understood (see BEHAVIOUR, ANIMAL: *Dormancy*).

Inactivity in response to adverse summer conditions (heat, drought, lack of food) is termed estivation. Estivation in some species is simply prolonged rest, usually in a favourable microhabitat; other estivating mammals regulate the metabolism, although the effects are typically not so pronounced as in hibernation.

Behavioral response to adverse conditions may involve the selection or construction of a suitable microhabitat (such as the cool, moist burrows of desert rodents). Migration is a second kind of behavioral response. The most obvious kind of mammalian migration is latitudinal. Many temperate-zone bats, for example, undertake extensive migrations, although other bat species hibernate near the summer foraging grounds in caves or other equable shelters during severe weather when insects are not available. Caribou (*Rangifer tarandus*) migrate from the tundra to the forest edge in search of suitable winter range, and a number of cetaceans and pinnipeds undertake long migrations from polar waters to more temperate latitudes. Gray whales (*Eschrichtius robustus*), for example, migrate southward to calving grounds along the coasts of South Korea and Baja California from summer feeding grounds in the Okhotsk, Bering, and Chukchi seas. Of comparable extent is the dispersive feeding migration of the northern fur seal (*Callorhinus ursinus*).

Migrations of lesser extent include the altitudinal movements of some ungulates, the American elk or wapiti (*Cervus canadensis*) and bighorn sheep (*Ovis canadensis*), for example, and the local migrations of certain bats from summer roosts to hibernacula. Most migratory patterns of mammals are part of a recurrent annual cycle, but the irruptive emigrations of lemmings and snowshoe hares are largely acyclic responses to population pressure on food supplies.

Populations. A population consists of individuals of three "ecological ages"—prereproductive, reproductive, and postreproductive. The structure and dynamics of a population depend, among other things, on the relative lengths of these ages, the rate of recruitment of individuals (either by birth or by immigration), and the rate of emigration or death. The reproductive potential of some rodents (particularly muroids) is well known; some mice are reproductively mature at four weeks of age, have gestation periods of three weeks or less, and may experience postpartum estrus, with the result that pregnancy and lactation may overlap. Litter size, moreover, may average four or more, and breeding may occur throughout the year in favourable localities. The reproductive potential of a species is, of course, a theoretical maximum that is rarely met inasmuch as, among other reasons, a given female typically does not reproduce throughout the year. Growth of a population depends on the survival of individuals to reproductive age. The absolute age at sexual maturity ranges from less than four weeks in some rodents to some 15 years in the African elephant (*Loxodonta africana*).

Postreproductive individuals are rare in most mammalian populations. Survival through more than a single reproductive season is probably uncommon in many small kinds, such as mice and shrews. Larger species typically have longer life spans than do smaller kinds, but some bats are known, on the basis of banding records, to live nearly 20 years. Many species show greater longevity in captivity than in the wild. Captive spiny anteaters (*Tachyglossus*) are reported to have lived more than 50 years. Horses have been reported to live more than 60 years, and the oldest known elephant lived to be 81. Man has greater potential longevity than any other known species.

Locomotion. Specialization in habitat preference has been accompanied by locomotor adaptations. Convergent evolution within a given adaptive mode has contributed to the ecological similarity of regional mammalian faunas. Terrestrial mammals have a number of modes of progression. The primitive mammalian stock was doubtless ambulatory and plantigrade, walking with the digits, metacarpals, and metatarsals (bones of the midfoot), and parts of the ankle and wrist in contact with the ground. The limbs of ambulatory mammals are typically mobile, capable of considerable rotation.

Mammals modified for running are termed cursorial. The stance of cursorial species may be digitigrade (the complete digits contacting the ground, as in dogs) or unguligrade (only tips of digits contacting the ground, as in horses). In advanced groups, limb movement consists of a single forward and backward direction.

Saltatory (leaping) locomotion, sometimes called "ricochetal," has arisen in several unrelated groups and typically is found in mammals of open habitats. Jumping mammals (some marsupials, lagomorphs, and several independent lineages of rodents) typically have elongate, plantigrade hind feet, reduced forelimbs, and long tails.

Mammals of several orders have attained great size and have converged on specializations for "graviportal" (ponderous) locomotion. Such kinds have no digit reduction and deploy the digits in a circle around the axis of the limb for maximum support, like the pedestal of a column.

The bats are the only truly flying mammals. Only with active flight have the resources of the aerial habitat been successfully exploited. Mammals of several kinds (dermopterans, marsupials, rodents) are adapted for gliding. A gliding habit is frequently accompanied by scansorial (climbing) locomotion. Many nongliders, such as tree squirrels, are also scansorial.

Well-adapted arboreal mammals frequently are plantigrade, five-toed, and equipped with highly mobile limbs. Many New World monkeys have a prehensile tail, which is used as a "fifth hand." Gibbons have reduced the primitive, opposable anthropoid thumb as a specialization for brachiation, or "arm-walking," in which the animal hangs from branches and moves by a series of long swings. The highly arboreal tarsiers have expanded pads on the digits to improve the grasp, and many other arboreal kinds have claws (sloths) or well-developed nails.

Several groups of mammals have independently assumed aquatic habits. In some cases, semi-aquatic or aquatic mammals are relatively unmodified representatives of otherwise terrestrial groups (for example, otters, muskrats, water shrews). Other kinds have undergone profound modification for natatorial locomotion and a pelagic habit. Pinniped carnivores (walruses and seals) give birth to their young on land, but cetaceans are completely helpless out of water, on which they depend for mechanical support and thermal insulation.

Food habits. The earliest mammals, like their reptilian forebears, were active predators. From such a basal stock there has been a complex radiation of trophic adaptations. Modern mammals occupy a wide spectrum of feeding niches. In most terrestrial and some aquatic communities, carnivorous mammals are at the top of the food pyramid. There are also mammalian primary consumers in most ecosystems. The voracious shrews, smallest of mammals, sometimes prey on vertebrates larger than themselves. They may eat twice their weight in food each day to maintain their active metabolism and compensate for heat loss caused by an unfavourable surface-to-volume ratio. The largest of vertebrates, the blue whale, feeds on krill, minute planktonic crustaceans.

Within a given lineage, the adaptive radiation in food habits may be broad. Some of the Carnivora have become omnivorous (raccoons, bears) or largely herbivorous (giant panda). Marsupials exhibit a great variety of feeding types, and in Australia marsupials have radiated to fill ecological niches highly analogous to those of placental mammals elsewhere; there are marsupial "moles," "anteaters," "mice," "rats," "cats," and "wolves." Some bandicoots have ecological roles similar to those of rabbits, and wombats are semifossorial (*i.e.*, partially burrowing) herbivores

Migrations

Aquatic habits

FORM AND FUNCTION

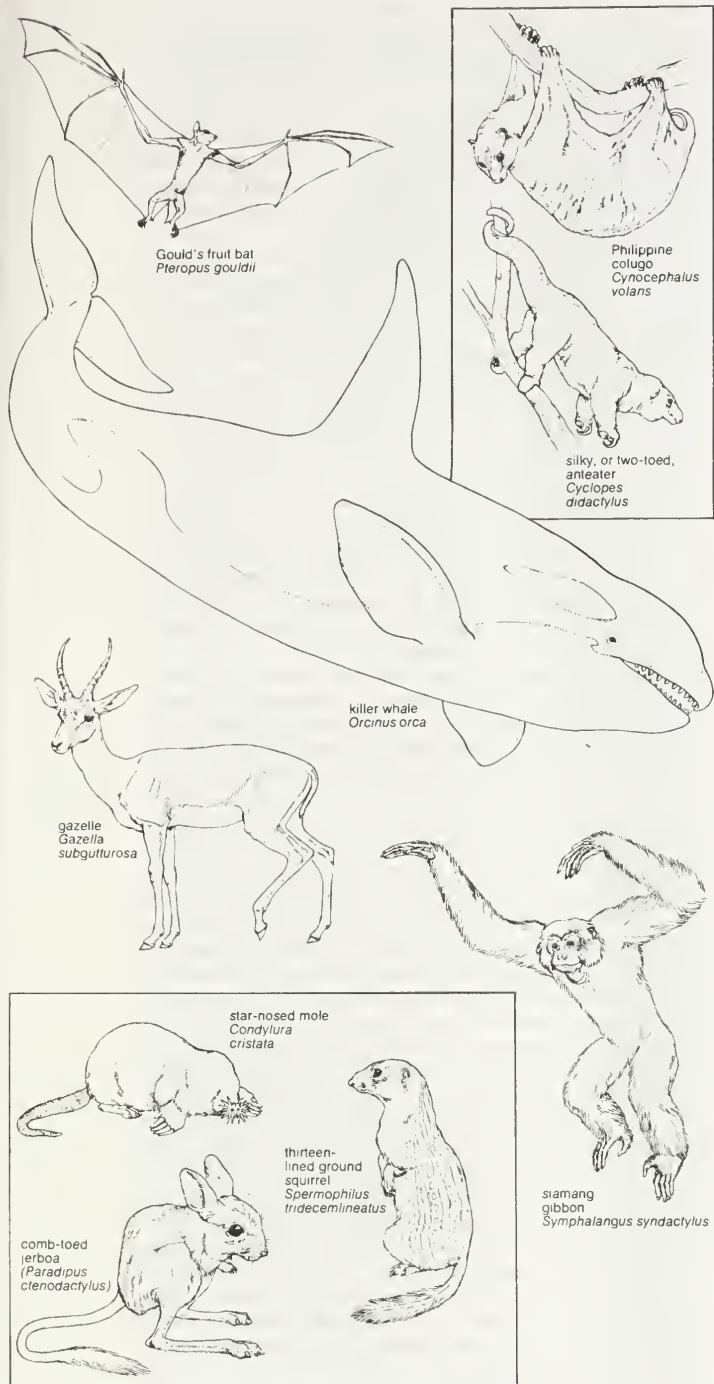


Figure 1: Range of body plans in mammals.

Drawing by R. Keane based on photographs courtesy of (Gould's fruit bat) J. Warham, (Philippine colugo) J.N. Hamlet, (silky, or two-toed anteater) New York Zoological Society, (comb-toed jerboa) Zeitschrift für Säugetierkunde

analogous to marmots. In Australia the niche of large grazers and browsers is filled by a variety of kangaroos and wallabies.

Radiation of bats

Within the bats there has also been a remarkable adaptive radiation in food habits. Early in the history of the order there evidently was a divergence into insectivorous and frugivorous lines. The frugivorous line (Megachiroptera) has generally maintained its fruit-eating habit, although some have become rather specialized nectar-feeders. The Microchiroptera have been less conservative and have undergone considerable divergence in feeding habits. A majority of living microchiropterans are insectivorous, but members of two different families have become fish eaters. Within the large neotropical family Phyllostomidae there are groups specialized to feed on fruit, nectar, insects, and small vertebrates (including other bats). A aberrant members of the family are the vampire bats (desmodontines) with a specialized dentition to aid the blood-lapping habit.

The skin and hair. The skin of mammals is constructed of two layers, a superficial nonvascular epidermis and an inner layer, the dermis or corium. The two layers interdigitate in dermal papillae (fingerlike projections), ridges of sensitive vascular dermis projecting into the epidermis. The outermost layers of the epidermis are cornified (i.e., impregnated with various tough proteins) and enucleate (lacking cell nuclei). The epidermis is composed of flattened (squamous) cells in layers and is the interface between the individual and the environment. Its primary function is defensive, and it is cornified to resist abrasion. The surface of the skin is coated with lipids and organic salts, the so-called "acid mantle," thought to have anti-fungal and antibacterial properties. Deep in the epidermis is an electronegative layer, a further deterrent to foreign organic or ionic agents.

The dermis lies beneath the epidermis and nourishes it. The cutaneous circulation of the dermis is variously developed in mammals but it is typically extensive, out of proportion to the nutritional needs of the tissue. The major role of the cutaneous circulation is to moderate body temperature and blood pressure by forming a peripheral shunt, an alternate route for the blood. Also in the dermis are sensory nerve endings to alert the individual to pressure (touch), heat, cold, and pain. In general, skin bearing hairs has few or no specialized sensory endings, but hairless skin, such as the lips and fingertips of man, has specialized endings. The sensation of touch on hairy skin in man depends on stimulation of the nerve fibres of the hairs.

Hair is derived from an invagination (pocketing) of the epidermis termed a follicle. Collectively, the hair is called the pelage. The individual hair is a rod of keratinized cells that may be cylindrical or more-or-less flattened. Keratin is a protein also found in claws and nails. The inner medulla of the hair is hollow and contains air; in the outer cortex layer there are frequently pigment granules. Associated with the hair follicle are nerve endings and a muscle, the arrector pili. The latter allows the erection of individual hairs to alter the insulative qualities of the pelage. The follicle also gives rise to sebaceous glands that produce sebum, a substance that lubricates the hair.

Most mammals have three distinct kinds of hairs. Guard hairs protect the rest of the pelage from abrasion and frequently from moisture, and usually lend a characteristic colour pattern. The thicker underfur is primarily insulative and may differ in colour from the guard hairs. The third common hair-type is the vibrissa, or whisker, a stiff, typically elongate, hair that functions in tactile sensation. Hairs may be further modified to form rigid quills. The "horn" of the rhinoceros is composed of a fibrous keratin material derived from hair. Examples of keratinized derivatives of the integument other than hair are horns, hooves, nails, claws, and baleen.

Types of hair

Even though the primary function of the skin is defensive, it has been modified in mammals to serve such diverse functions as thermoregulation and nourishment of young. Secretions of sweat glands promote cooling due to evaporation at the surface of the body, and mammary glands are thought to be derived from sweat glands.

In certain groups (primates in particular) the skin of the face is under intricate muscular control and movements of the skin express and communicate "emotion." In many mammals the colour and pattern of the pelage is important in communicative behaviour. Patterns may be dymanitic (startling, such as the mane of the male lion or hamadryas baboon), sematic (warning, such as the bold pattern of skunks), or cryptic (concealing), perhaps the most common adaptation of pelage colour.

Hair has been secondarily lost or considerably reduced in some kinds of mammals. In adult cetaceans insulation is provided by thick subcutaneous fat deposits, or blubber, with hair limited to a few stiff vibrissae about the mouth. The naked skin of whales is one of a number of features that contribute to the remarkably advanced hydrodynamics of locomotion in the group. Some fossorial mammals also tend toward reduction of the hair. This is shown most strikingly by the sand rats (*Heterocephalus*) of northeast

Africa, but considerable loss of hair has also occurred in some species of pocket gophers (Geomyidae). Hair may also be lost on restricted areas of the skin, as from the face in many monkeys, or the buttocks of mandrills, and may be sparse on elephants and such highly modified kinds as pangolins and armadillos.

Indeterminate (continuous) growth, as seen in the hair of the head in man, is rare among mammals. Hairs with determinate growth are subject to wear and must be replaced periodically—a process termed molt. The first coat of a young mammal is referred to as the juvenal pelage, which typically is of fine texture like the underfur of adults and is replaced by a postjuvenal molt. Juvenal pelage is succeeded by the subadult pelage, which in some species is not markedly distinct from that of the adult, or directly by adult pelage. Once this pelage is acquired molting continues to recur at intervals, often annually or semi-annually, sometimes more frequently. The pattern of molt typically is orderly, but varies widely between species. Some mammals apparently molt continuously, with a few hairs replaced at a time throughout the year.

Dentition. Specialization in food habits has led to profound dental changes. The primitive mammalian tooth had high, sharp cusps and served to tear flesh or crush chitinous material (primarily the exoskeletons of terrestrial arthropods, such as insects). Herbivores tend to have specialized cheek teeth with complex occlusal (contact) patterns and various ways of expanding the crown and circumventing the problem of wear. Omnivorous mammals, such as bears, pigs, and man, tend to have molars with low, rounded cusps, termed "bunodont."

Anteating

A prime example of convergence in conjunction with dietary specialization is seen in those mammals adapted to feeding on social insects (generally termed myrmecophagy, "anteating"). This habit has led to remarkably similar morphology in such diverse groups as the monotreme echidnas (Tachyglossidae), one dasyurid marsupial (Myrmecobinae), some edentates (Myrmecophagidae) the armadillo, and pangolins (Pholidota). Trends frequently associated with myrmecophagy include: strong claws, an elongate, terete skull, a vermiform, extensible tongue, marked reduction in the mandible, and loss or extreme simplification of the dentition.

Specialized herbivores evolved early in mammalian history. The earliest (and with the longest evolutionary history) were the multituberculates. Evolutionary convergence in teeth (*i.e.*, resemblances not due to common ancestry) has occurred widely in herbivorous groups; most have incisors modified for nipping or gnawing, have lost teeth with the resultant development of a gap (diastema) in the tooth row, and exhibit some molarization (expansion and flattening) of premolars to expand the grinding surface of the cheek teeth. Rootless incisors or cheek teeth have evolved frequently, an open pulp cavity allowing continual growth throughout life. Herbivorous specializations have evolved independently in multituberculates, rodents, lagomorphs, primates, and in the wide radiation of ungulate and subungulate orders.

Skeleton. The mammalian skeletal system shows a number of advances over that of lower vertebrates. The mode of ossification of the long bones is characteristic. In lower vertebrates each long bone has a single centre of ossification, the diaphysis, and replacement of cartilage by bone proceeds from the centre toward the ends, which may remain cartilaginous, even in adults. In mammals secondary centres of ossification, the epiphyses, develop at the ends of the bones. Growth of bones occurs in zones of cartilage between diaphysis and epiphyses. Mammalian skeletal growth is termed determinate, for once the actively growing zone of cartilage is obliterated, growth in length ceases. As in all bony vertebrates, of course, there is continual renewal of bone throughout life. The advantage of epiphyseal ossification lies in the fact that the bones have strong articular surfaces before the skeleton is mature. In general, the skeleton of the adult mammal has less structural cartilage than does that of a reptile.

The skeletal system of mammals and other vertebrates is broadly divisible functionally into axial and appendicular portions. The axial skeleton consists of the braincase

(cranium) and the backbone and ribs, and serves primarily to protect the central nervous system. The limbs and their girdles constitute the appendicular skeleton. In addition, there are skeletal elements derived from the gill arches of primitive vertebrates, collectively termed the visceral skeleton. Visceral elements in the mammalian skeleton include the jaws, the hyoid apparatus supporting the tongue, and the auditory ossicles of the middle ear. The postcranial axial skeleton in mammals generally has remained rather conservative during the course of evolution. The vast majority of mammals have seven cervical (neck) vertebrae; exceptions are sloths, with six or nine cervicals, and the Sirenia with six. The anterior two cervical vertebrae are differentiated as atlas and axis. Specialized articulations of these two bones allow complex movements of the head on the trunk. Thoracic vertebrae bear ribs and are variable in number. The anterior ribs converge toward the ventral midline to articulate with the sternum, or breastbone, forming a semirigid thoracic "basket" for the protection of heart and lungs. Posterior to the thoracic region are the lumbar vertebrae, ranging from two to 21 in number (most frequently four to seven). Mammals have no lumbar ribs. There are usually three to five sacral vertebrae, but some edentates have as many as 13. Sacral vertebrae fuse to form the sacrum, to which the pelvic girdle is attached. Caudal (tail) vertebrae range in number from five (fused elements of the coccyx of man) to 50.

The basic structure of the vertebral column is comparable throughout the Mammalia, although in many instances modifications have occurred in specialized locomotor modes to gain particular mechanical advantages. The vertebral column and associated muscles of many mammals is structurally analogous to a cantilever girder.

The skull is composite in origin and complex in function. Functionally the bones of the head are separable into the braincase and the jaws. In general, it is the head of the animal that meets the environment. The skull protects the brain and sense capsules, houses the teeth and tongue, and the entrance to the pharynx. Thus the head functions in sensory reception, food acquisition, defense, respiration, and (in higher groups) communication. To serve these functions, bony elements have been recruited from the visceral skeleton, the endochondral skeleton, and from the dermal skeleton of lower vertebrates.

The skull

The skull of mammals differs markedly from that of reptiles because of the great expansion of the brain. The sphenoid bones that form the reptilian braincase form only the floor of the braincase in mammals. The side is formed in part by the alisphenoid bone, derived from the epipterygoid, a part of the reptilian palate. Dermal elements, the frontals and parietals, have come to lie deep to (beneath) the muscles of the jaw to form the dorsum of the braincase. Reptilian dermal roofing bones, lying superficial to the muscles of the jaw, are represented in mammals only by the jugal bone of the zygomatic arch, which lies under the eye.

In mammals, a secondary palate is formed by processes of the maxillary bones and the palatines, with the pterygoid bones reduced in importance. The secondary palate separates the nasal passages from the oral cavity and allows continuous breathing while chewing or suckling.

Other specializations of the mammalian skull include paired occipital condyles (the articulating surfaces at the neck) and an expanded nasal chamber with complexly folded turbinal bones, providing a large area for detection of odours. Eutherians have evolved bony protection for the middle ear, the auditory bulla. The development of this structure varies, although an annular (ring-shaped) tympanic bone is always present.

The bones of the mammalian middle ear are a diagnostic feature of the class. The three auditory ossicles form a series of levers that serve mechanically to increase the amplitude of sound waves reaching the tympanic membrane, or eardrum, produced as disturbances of the air. The innermost bone is the stapes, or "stirrup bone." It rests against the oval window of the inner ear. The stapes is homologous with the entire stapedial structure of reptiles, which, in turn, was derived from the hyomandibular arch of primitive vertebrates. The incus, or "anvil," ar-

Bones of the ear

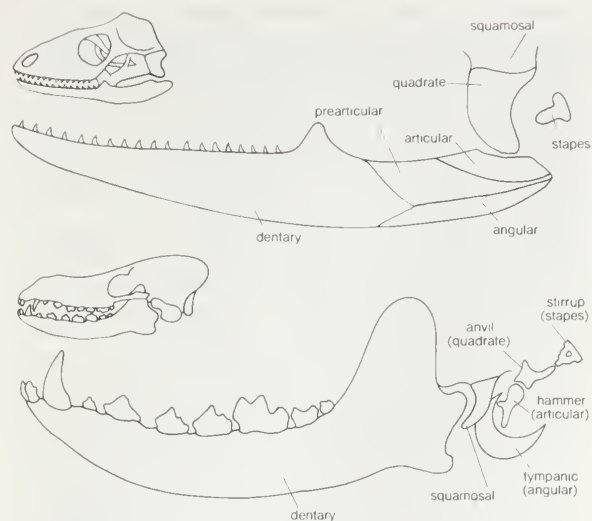


Figure 2: Comparison of lower jaw and ear region in the skull of (top) a reptile and (bottom) a mammal.

articulates with the stapes. The incus was derived from the quadrate bone, which is involved in the jaw articulation in reptiles. The malleus, or "hammer," rests against the tympanic membrane and articulates with the incus. The malleus is the homologue of the reptilian articular bone. The mechanical efficiency of the middle ear has thus been increased by the incorporation of two bones of the reptilian jaw assemblage. In mammals the lower jaw is a single bone, the dentary, which articulates with the squamosal of the skull.

The limbs and girdles have been greatly modified with locomotor adaptations. The ancestral mammal had well-developed limbs and was five-toed. In each limb there were two distal (outer) elements (radius and ulna in the forelimb; tibia and fibula in the hindlimb) and a single proximal (inner or upper) element (humerus; femur). There were nine bones in the wrist, the carpals; and seven bones in the ankle, the tarsals. The phalangeal formula (the number of phalangeal bones in each digit, numbered from inside outward) is 2-3-3-3-3 in primitive mammals; in primitive reptiles it is 2-3-4-5-3. Modifications in mammalian limbs have involved reduction, loss, or fusion of bones. Loss of the clavicle from the shoulder girdle, reduction in the number of toes, and modifications of tarsal and carpal bones are typical correlates of cursorial locomotion. Scansorial and arboreal groups tend to maintain or emphasize the primitive divergence of the thumb and hallux (the inner toe on the hindfoot). Some details of limb modifications are included in the synopsis of the orders of mammals (see below *Annotated classification*).

Centres of ossification sometimes develop in nonbony connective tissue. Such bones are termed heterotopic or sesamoid elements. The kneecap (patella) is such a bone. Another important bone of this sort, found in many kinds of mammals, is the baculum, or os penis, which occurs as a stiffening rod in the penis of such groups as carnivores, many bats, rodents, some insectivores, and many primates. The os clitoridis is a homologous structure found in females.

Muscles. The muscular system of mammals is generally comparable to that of reptiles. With changes in locomotion, the proportions and specific functions of muscular elements have been altered, but the relationships of these muscles remain essentially the same. Exceptions to this generalization are the muscles of the skin and of the jaw.

The panniculus carnosus is a sheath of dermal (skin) muscle, developed in many mammals, which allows the movement of the skin independent of the movement of deeper muscle masses. Such movement functions in such mundane activities as the twitching of the skin to foil insect pests and in some species also is important in shivering, a characteristic heat-producing response to thermal stress. The dermal musculature of the facial region is particularly well developed in primates and carnivores, but occurs in other groups as well. Facial mobility allows expression

that may be of importance in the behavioral maintenance of interspecific social structure.

The temporalis muscle is the major adductor (closer) of the reptilian jaw. In mammals the temporalis is divided into a deep temporalis proper and a more superficial masseter muscle. The temporalis attaches to the coronoid process of the mandible (lower jaw) and the temporal bone of the skull. The masseter passes from the angular process of the mandible to the zygomatic arch. The masseter allows an anteroposterior (forward-backward) movement of the jaw and is highly developed in mammals, such as rodents, in which grinding is the important function of the dentition.

Viscera. Digestive system. The alimentary canal is highly specialized in many kinds of mammals. In general, specializations of the gut accompany herbivorous habits. The intestines of herbivores are typically elongate and the stomach may also be specialized. Subdivision of the gut allows areas of differing physiological environments for the activities of different sorts of enzymes and symbiotic bacteria (which aid the animal by breaking down certain "indigestible" compounds). In ruminant artiodactyls the stomach has up to four chambers, each with a particular function in the processing of vegetable material. A cecum is common in many herbivores. The cecum is a blind sac at the distal end of the small intestine where complex compounds like cellulose are acted upon by symbiotic bacteria. The vermiform appendix is a diverticulum of the cecum. The appendix is rich in lymphoid tissue and in many mammals is concerned with defense against toxic bacterial products.

Hares and rabbits, the sewellel or "mountain beaver" (*Aplodontia rufa*), and some insectivores exhibit a phenomenon, known as reingestion, in which at intervals specialized fecal pellets are produced, which are taken in the mouth and passed through the alimentary canal a second time. Where known to be present, this pattern seems to be obligatory. Reingestion is poorly understood, but it is generally supposed that the process allows the animal to absorb in the upper gut vitamins produced by the microflora of the lower gut but not absorbable there.

Excretory system. The mammalian kidney is constructed of a large number of functional units called nephrons. Each nephron consists of a distal tubule, a medial section termed the loop of Henle, a proximal tubule, and a renal corpuscle. The renal corpuscle is a knot of capillaries, called a glomerulus, surrounded by a sheath, Bowman's capsule. The renal corpuscle is a pressure filter, relying on blood pressure to remove water, ions, and small organic molecules from the blood. Some of the material removed is waste, but some is of value to the organism. The filtrate is sorted by the tubules, and water and needed solutes are reabsorbed. Reabsorption is both passive (osmotic) and active (based on ion transport systems). The distal convoluted tubules drain into collecting tubules which, in turn, empty into the calyces, or branches, of the renal pelvis, the expanded end of the ureter. The pressure-pump nephron of mammals is so efficient that the renal portal system of lower vertebrates has been completely lost. Mammalian kidneys show considerable variety in structure, relative to the environmental demands on a given species. In particular, desert rodents have long loops of Henle and are able to reabsorb much water and to excrete a highly concentrated urine. Urea is the end product of protein metabolism in mammals, and excretion is therefore called urcotelic.

Reproductive system. The testes of mammals descend from the abdominal cavity to lie in a compartmented pouch termed the scrotum. In some species the testes are permanently scrotal and the scrotum is sealed off from the general abdominal cavity. In others the testes migrate to the scrotum only during the breeding season. It is thought that the temperature of the abdominal cavity is too high to allow spermatogenesis; the scrotum allows cooling of the testes.

The transport of spermatozoa is comparable to that in reptiles, relying on ducts derived from urinary ducts of earlier vertebrates. Mammalian specialities are the bulbourethral (or Cowper's) glands, the prostate gland, and

Muscles
of the jaw

Action
of the
kidney

the seminal vesicle or vesicular gland. Each of these glands adds secretions to the spermatozoa to form semen, which passes from the body via a canal (urethra) in the highly vascular, erectile penis. The tip of the penis, the glans, may have a complex morphology and has been used as a taxonomic character in some groups. The penis may be retracted into a sheath along the abdomen or may be pendulous, as in bats and many primates.

Types of uterus

The structure of the female reproductive tract is variable. Four uterine types are generally recognized among placentals, based on the relationship of the uterine horns (branches). A duplex uterus characterizes rodents and rabbits; the uterine horns are completely separated and have separate cervixes opening into the vagina. Carnivores have

From (A) A.S. Romer *The Vertebrate Body*, 4th ed (1970) W.B. Saunders Company Philadelphia (B, C) W.F. Walker *Vertebrate Dissection*, 4th ed (1970) W.B. Saunders Company Philadelphia

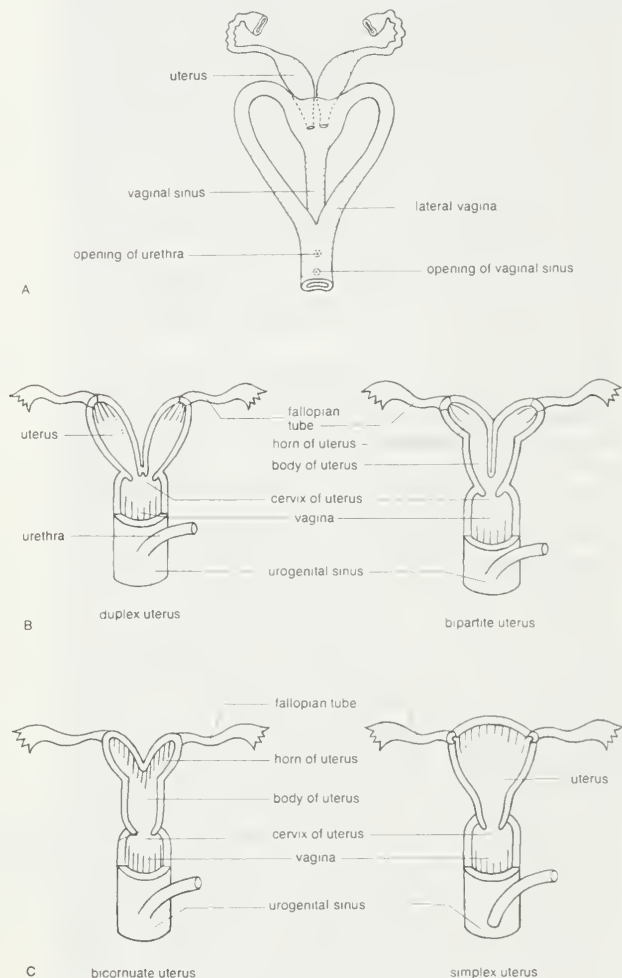


Figure 3. Female reproductive systems. (A) Duplex uterus with median vaginal tube (kangaroo). (B, C) Placental uteri.

a bipartite uterus, in which the horns are largely separate but enter the vagina by a single cervix. In the bicornuate uterus, typical of many ungulates, the horns are distinct for less than half their length; the lower part of the uterus is a common chamber, the body. Higher primates have a simplex uterus in which all separation between the horns is lacking, forming a single chamber.

The female reproductive tract of marsupials is termed didelphous; the vagina is paired, as are oviducts and uteri. In primitive marsupials there are paired vaginae lateral to the ureters. In more advanced groups, such as kangaroos, the lateral vaginae persist and conduct the migration of spermatozoa, but a medial "pseudovagina" functions as the birth canal.

Monotremes have paired uteri and oviducts, which empty into a urogenital sinus (cavity) as do fluid wastes. The sinus passes into the cloaca, a common receptacle for reproductive and excretory products.

Circulatory system. In mammals as in birds, right and

left ventricles of the heart are completely separated, so that pulmonary (lung) and systemic (body) circulations are completely independent. Oxygenated blood arrives in the left atrium from the lungs and passes to the left ventricle, whence it is forced through the aorta to the systemic circulation. Deoxygenated blood from the tissues returns to the right atrium via a large vein, the vena cava, and is pumped to the pulmonary capillary bed through the pulmonary artery.

Among vertebrates, contraction of the heart is myogenic; rhythm is inherent in all cardiac muscle, but in myogenic hearts the pacemaker is derived from cardiac tissue. The pacemaker in mammals (and also in birds) is an oblong mass of specialized cells called the sinoatrial node, located in the right atrium near the junction with the venae cavae. A wave of excitation spreads from this node to the atrioventricular node, which is located in the right atrium near the base of the interatrial septum. From this point excitation is conducted along the atrioventricular bundle (bundle of His) and enters the main mass of cardiac tissue along fine branches, the Purkinje fibres. Homeostatic control of the heart by neuroendocrine or other agents is mediated through the intrinsic control network of the heart.

Blood leaves the left ventricle through the aorta. The mammalian aorta is an unpaired structure derived from the left fourth aortic arch of the primitive vertebrate. Birds, on the other hand, retain the right fourth arch.

The circulatory system forms a complex communication and distribution network to all physiologically active tissues of the body. A constant, copious supply of oxygen is required to sustain the active, endothermic physiology of the higher vertebrates. The efficiency of the four-chambered heart is important to this function. Oxygen is transported by specialized red blood cells, or erythrocytes, as in all vertebrates. Packaging the oxygen-bearing pigment hemoglobin in erythrocytes keeps the viscosity of the blood minimal and thereby allows efficient circulation while limiting the mechanical load on the heart. The mammalian erythrocyte is a highly evolved structure; its discoid, biconcave shape allows maximal surface area per unit volume. When mature and functional, mammalian red blood cells are enucleate.

Respiratory system. Closely coupled with the circulatory system is the ventilatory (breathing) apparatus, the lungs and associated structures. Ventilation in mammals is unique. The lungs themselves are less efficient than those of birds, for air movement consists of an ebb and flow, rather than a one-way circuit, so a residuum of air remains that cannot be expired. Ventilation in mammals is by means of a negative pressure pump made possible by the development of a definitive thoracic cavity with the evolution of the diaphragm.

The diaphragm is a unique, composite structure consisting of (1) the transverse septum (a wall that primitively separates the heart from the general viscera); (2) pleuroperitoneal folds from the body wall; (3) mesenteric folds; and (4) axial muscles inserting on a central tendon, or diaphragmatic aponeurosis.

The lungs lie in separate, airtight compartments called pleural cavities, separated by the mediastinum. As the size of the pleural cavity is increased the lung is expanded and air flows in passively. Enlargement of the pleural cavity is produced by contraction of the diaphragm or by elevation of the ribs. The relaxed diaphragm domes upward, but when contracted it stretches flat. Expiration is an active movement brought about by contraction of abdominal muscles against the viscera.

Air typically enters the respiratory passages through the nostrils where it may be warmed and moistened. It passes above the bony palate and the soft palate and enters the pharynx. In the pharynx the passages for air and food cross. Air enters the trachea, which divides at the level of the lungs into primary bronchi. A characteristic feature of the trachea of many mammals is the larynx. Vocal cords stretch across the larynx and are vibrated by forced expiration to produce sound. The laryngeal apparatus may be greatly modified for the production of complex vocalizations. In some groups, for example, howler monkeys

Characteristics of the blood

(*Alouatta*), the hyoid apparatus is incorporated into the sound-producing organ, as a resonating chamber.

Nervous and endocrine systems. The nervous system and the system of endocrine glands are closely related to one another in their function, for both serve to coordinate activity. The endocrine glands of mammals generally have more complex regulatory functions than do those of lower vertebrates. This is particularly true of the pituitary gland, which supplies hormones that regulate the reproductive cycle. Follicle stimulating hormone (FSH) initiates the maturation of the ovarian follicle. Luteinizing hormone (LH) mediates the formation of the corpus luteum from the follicle following ovulation. Prolactin, also a product of the anterior pituitary, stimulates the secretion of milk.

Control of the pituitary glands is partially by means of neurohumours from the hypothalamus, a part of the fore-brain in contact with the pituitary gland by nervous and circulatory pathways. The hypothalamus is of the utmost importance in mammals, for it integrates stimuli from both internal and external environments, channelling signals to higher centres or into autonomic pathways.

The cerebellum of vertebrates is at the anterior end of the hindbrain. Its function is to coordinate motor activities and to maintain posture. In most mammals the cerebellum is highly developed and its surface may be convoluted to increase its area. The data with which the cerebellum works arrive from proprioceptors ("self-sensors") in the muscles and from the membranous labyrinth of the inner ear, the latter giving information on position and movements of the head.

In the vertebrate ancestors of the mammals the cerebral hemispheres were centres for the reception of olfactory stimuli. Vertebrate evolution has favoured an increasing importance of these lobes in the integration of stimuli. Their great development in mammals as centres of association is responsible for the "creative" behaviour of members of the class; *i.e.*, the ability to learn, to adapt as individuals to short-term environmental change through appropriate responses on the basis of previous experience. In vertebrate evolution the gray matter of the cerebrum has moved from a primitive internal position in the hemispheres to a superficial position. The superficial gray matter is termed the pallium. The paleopallium of amphibians has become the olfactory lobes of the higher vertebrates; the dorsolateral surface, or archipallium, the mammalian hippocampus. The great neural advance of the mammals lies in the elaboration of the neopallium, which makes up the bulk of the cerebrum. The neopallium is an association centre, the dominant centre of neural function, and is involved in so-called "intelligent" response. By contrast, the highest centre in the avian brain is the corpus striatum, an evolutionary product of the basal nuclei of the amphibian brain. The bulk of the complex behaviour of birds is instinctive. The surface of the neopallium tends in some mammals to be greatly expanded by convoluting, forming folds (gyri) between deep grooves (sulci).

Evolution and classification

THE EVOLUTION OF THE MAMMALIAN CONDITION

Mammals were derived in the Triassic Period from members of the reptilian order Therapsida. The therapsids, members of the subclass Synapsida (sometimes called the mammal-like reptiles), generally were unimpressive in relation to other reptiles of their time. Synapsids were present in the Carboniferous Period (about 280,000,000 to 345,000,000 years ago) and are one of the earliest known reptilian groups. They were the dominant reptiles of Permian times (about 225,000,000 to 280,000,000 years ago), and although they were primarily predacious in habit, the adaptive radiation included herbivorous species as well. In the Mesozoic (about 225,000,000 to 65,000,000 years ago), the importance of the synapsids was generally assumed by the archosaurs or "ruling reptiles," the therapsids, in general, being small, active carnivores. Therapsids tended to evolve a specialized heterodont dentition and to improve the mechanics of locomotion by bringing the plane of action of the limbs close to the trunk. A secondary palate was developed and the temporal musculature was expanded.

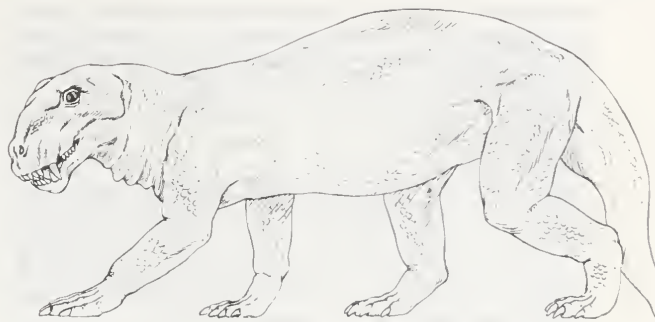


Figure 4: Reconstruction of *Lycaenops*, a Late Permian mammal-like reptile.

Drawing by R. Keane based on A.S. Romer, *The Vertebrate Story*, copyright 1959, the University of Chicago Press

The several features that separate modern reptiles from modern mammals doubtless evolved at different rates. Many attributes of mammals are correlated with their highly active habit; for example, efficient double circulation with a completely four-chambered heart, anucleate and biconcave erythrocytes, the diaphragm, and the secondary palate (which separates passages for food and air and allows breathing during mastication or suckling). Hair for insulation is a correlate of endothermy, the physiological maintenance of individual temperature independent of environmental temperature, and endothermy allows high levels of sustained activity. The unique characteristics of mammals thus would seem to have evolved as a complex interrelated system.

Because the characteristics that separate reptiles and mammals evolved at different rates and in response to a variety of interrelated conditions, at any point in the period of transition from reptiles to mammals there were forms that combined various characteristics of both groups. Such a pattern of evolution is termed "mosaic" and is a common phenomenon in those transitions marking the origin of major new adaptive types. To simplify definitions and to allow the strict delimitation of the Mammalia, some authors have suggested basing the boundary on a single character, the articulation of the jaw between the dentary and squamosal bones and the attendant movement of accessory jaw bones to the middle ear as auditory ossicles. The use of a single osteological character allows the placement in a logical classification of numerous fossil species, other mammalian characters of which, such as the degree of endothermy and nursing of young and the condition of the internal organs, probably never will be evaluated. It must be recognized, however, that were the advanced therapsids alive, taxonomists would be hard-put to decide which to place in the Reptilia and which in the Mammalia.

Mosaic evolution

CLASSIFICATION

Distinguishing taxonomic features. The higher classification of the Mammalia is based on consideration of a broad array of characters. Traditionally, evidence from comparative anatomy was of predominant importance, but more recently information from such disciplines as physiology, serology, and genetics has proved useful in considering relationships. Comparative study of living organisms is supplemented by the findings of palaeontology. Study of the fossil record adds a historical dimension to knowledge of mammalian relationships. In some cases, the horses for example, the fossil record has been adequate to allow lineages to be traced in great detail.

Relative to that of other major vertebrate groups, the fossil record of mammals is good. Fossilization depends upon a great many factors, most important among which are the structure of the organism, its habitat, and conditions at the time of death. The most common remains of mammals are teeth and the associated bones of the jaw and skull. Enamel covering the typical mammalian tooth is composed of prismatic rods of crystalline apatite and is the hardest tissue in the mammalian body. It is highly resistant to chemical and physical weathering. Because of the abundance of teeth in deposits of fossil mammals, dental characteristics have been stressed in the interpretation of mammalian phylogeny and relationships. Dental

features are particularly well suited for this important role in classification because they reflect the broad radiation of mammalian feeding specializations from the primitive predaceous habit.

Annotated classification. The following classification has wide acceptance.

CLASS MAMMALIA

Vertebrate animals with the lower jaw formed by a single bone (dentary), articulating directly with the squamosal bone of the skull; dentition generally heterodont, thecodont, and diphodont; body covered with hair derived from epidermis; young nourished by milk from mammary glands. For additional characteristics of the class, see *Form and function*, above. Groups preceded by a dagger mark (†) below are extinct.

Subclass Prototheria

Mammals in which principal cusps of molars form an antero-posterior row and much of the side of the braincase is formed by the periotic bone rather than by the alisphenoid. Primitive characters include uncoiled cochlea of the inner ear and absence of spine on the scapula (shoulder blade).

†*Infraclass Eotheria*

Fossil only; Upper Triassic to Lower Cretaceous of Eurasia and North America. Upper molars with internal and external cingula (shelves at base of crown), lower molars with internal cingulum only.

†*Order Triconodonta* (triconodonts). Fossil only; Northern Hemisphere; Upper Triassic to Lower Cretaceous. Probably carnivorous, the teeth well differentiated; molars typically with 3 prominent cusps; dentary bone lacking angular process.

†*Order Docodonta* (docodonts). Known only from fossil teeth and jaws from Upper Jurassic of Europe. Molars expanded, with transverse median furrow; angular process of mandible directed ventrally rather than posteriorly ("pseudangular").

Infraclass Ornithodelphia

Order Monotremata (monotremes). Pleistocene to Recent of Australia and New Guinea; possible record in Miocene of Australia family (Ectopodontidae). Reproduction by shelled eggs, incubated by female; no organized nipple, the milk flowing from mammary ducts onto the fur of the abdomen; cloaca present; adults toothless, rostrum covered with horny beak. Total length to 80 cm (31.5 in.). Examples: duck-billed platypus, and echidnas or spiny anteaters.

†*Infraclass Allotheria*

†*Order Multituberculata* (multituberculates). Fossil only; Upper Jurassic to lower Eocene of North America and Eurasia. The earliest herbivorous mammals; specializations of teeth and jaws superficially rodent-like; molars elongate, with two or three longitudinal rows of cusps; lower premolars tending to become serrated shearing blades. Size larger than typical Mesozoic mammals, to about that of a woodchuck.

Subclass Theria

†*Infraclass Trituberculata*

†*Order Symmetrodonta* (symmetrodonts). Fossil only; Upper Jurassic to Lower Cretaceous of Eurasia and North America. Molars with 3 prominent cusps arranged in symmetrical triangles; lower molars without posterior talonid basin ("heel").

†*Order Pantotheria* (pantotheres). Fossil only; Middle Jurassic to Lower Cretaceous of North America and Europe. Molars with 3 prominent anterior cusps (trigonid), asymmetrical; posterior basin (talonid) of lower molars developed, as in higher mammals.

Infraclass Metatheria

Order Marsupialia (marsupials). Upper Cretaceous to Recent of North and South America, Europe, Australia. Placenta formed from yolk sac (simple chorio-allantoic placenta in bandicoots), the young born in a relatively undeveloped state after brief gestation, developing further attached to mammae in marsupium (pouch) or along marsupial fold; female reproductive tract bifid, uterus and vagina double, penis bifurcate, posterior to serotum. Angular process of mandible inflected; auditory bullae (when present) formed from alisphenoid bone and thus not homologous with bullae of eutherians. Epipubis ("marsupial bone") present; length of head and body varying from about 10 cm (4 in.; small opossums) to more than 160 cm (63 in.; kangaroos). Examples: opossums, kangaroos, wallabies, bandicoots, wombats, phalangers, koala.

Infraclass Eutheria (placental mammals)

Order Insectivora (insectivores). Upper Cretaceous to Recent, worldwide except Australia, polar regions, and many oceanic islands. A diverse group of mammals sharing certain primitive characteristics but highly specialized in several adap-

tive modes: dentition generally heterodont, rooted, variously modified; molars usually tuberculosectorial. Stance plantigrade to semiplantigrade; habit fossorial, terrestrial, arboreal, semi-aquatic; insectivorous to omnivorous. Braincase low, rising but little above plane of rostrum; snout generally long, pointed, modified for olfaction and touch; zygomatic arch complete or incomplete. Length of head and body, 5 cm (2 in.; smallest shrews) to 45 cm (18 in.; largest hedgehogs). Shrews of the genera *Microsorex*, *Sorex*, and *Suncus* are the smallest of living mammals. Examples: shrews, moles, hedgehogs, tenrecs, solenodon.

Order Macroscelidea (elephant shrews). Lower Oligocene to Recent of Africa. Antermost incisors widely spaced (incisor series forming a straight, anteroposterior row); posterior incisors caniniform; upper canine premolariform, bicuspid; posterior premolars molariform, last premolar largest tooth in dental series; molars quadrate, reduced in size or absent posteriorly; auditory bullae well developed. Stance plantigrade, hind feet elongate, modified for hopping gait; tail elongate. Length of head and body of living species, 12 to 20 cm (5–8 in.); tail 12 to 18 cm (5–7 in.; treated in section on *Insectivora*).

Order Dermoptera (colugos or "flying lemurs"). Upper Paleocene to lower Eocene of North America (Plagiomenidae), Recent of Southeast Asia (Cynocephalidae). In living species, upper incisors reduced, lower incisors procumbent, spatulate, pectinate (comb-like), molars low-crowned, used for shearing the vegetable diet. The most specialized and efficient of mammalian gliders; the patagium (gliding membrane) extends from the neck to phalanges of hand, to phalanges of foot, to tip of tail; total length of head and body, 34–45 cm (13–18 in.).

†*Order Tillodontia* (tillodonts). Fossil only; upper Paleocene to middle Eocene of Northern Hemisphere. Omnivorous or herbivorous; incisors rootless, enlarged, molars low-crowned; perhaps superficially bearlike in appearance; total length to approximately 120 cm (47 in.).

†*Order Taeniodonta* (taeniodonts). Fossil only; lower Paleocene to upper Eocene of North America, middle Eocene of Asia. Large, browsing herbivores; outer incisors rodent-like, inner incisors lost; canines well developed; cheek teeth simple, rootless pegs; length of skull to about 35 cm (14 in.).

Order Chiroptera (bats). Middle Eocene to Recent in tropical and temperate regions; volant, the wings consisting of a fold of skin (patagium) supported by elongated bones of the second through fifth digits; thumb clawed in suborder Microchiroptera, thumb and index finger clawed in Megachiroptera; hind limbs variously specialized but generally weak. Incisors generally reduced, canines prominent, cheek teeth rather primitive, with prominent ectoloph. Length of head and body from about 3 to 40 cm (1.2–16 in.); wing span to about 120 cm (47 in.).

Order Primates (primates). Late Cretaceous to Recent of Eurasia, Africa, North and South America. A diverse group of generalized, generally omnivorous mammals, modified for arboreality (some secondarily terrestrial). Dentition little specialized; a pair of incisors lost above and below, canines prominent, sexually dimorphic in some species; premolars usually bicuspid, sometimes caniniform; molars typically low-crowned, tuberculosectorial to bunodont. Limbs and girdles adapted for flexibility with little loss or fusion of elements; clavicle present; digits 1 (hallux and pollex) divergent, tending to become opposable; stance plantigrade. Olfaction de-emphasized in favour of vision, and snout reduced; eyes rotated forward allowing trend toward stereoscopic vision; in more advanced groups cerebrum expanded to cover cerebellum, braincase expanded. Length of head and body (including extended hind limbs), 15 cm (6 in.; pygmy marmosets, tarsiers) to more than 200 cm (79 in.; *Homo*). Examples: lemurs, lorises, Old World and New World monkeys, marmosets, great apes, man.

Order Edentata (edentates). Lower Eocene to Recent of North and South America. Anterior teeth lost and cheek teeth reduced to simple rootless pegs without enamel, or (in anteaters) lost altogether. Secondary articulations (xenarthrous processes) present between vertebrae in addition to zygapophyses (the usual articulating surfaces); cervical vertebrae 6 to 9, in some cases fused; acromion and coracoid processes of scapula enlarged, often joining to form bony ring; claws generally stout. Several groups walk on the sides of feet and structure may be modified accordingly. Brain small, skull low, zygomatic arch usually incomplete. Total length of fossil edentates to 300 cm (118 in.); living species to 185 cm (73 in.; giant anteaters). Examples: sloths, armadillos, anteaters.

Order Pholidota (pangolins). Oligocene to Recent of Eurasia and Africa; possible North America representatives (Palaeodontidae) upper Paleocene to lower Oligocene. Body covered with heavy, overlapping scales with hairs only at bases of scales; tail prehensile in some species. Myrmecophagous (anteating); skull elongate, braincase low; lower jaw reduced to thin shaft

of bone; feet 5-toed but only 3 toes of forefeet strongly clawed. Total length 75 to 170 cm (30–67 in.).

Order Tubulidentata (aardvark). Miocene to Recent of Africa, Pliocene of Eurasia; possible representatives in Eocene and Oligocene of Europe. Teeth reduced to 4 or 5 simple pegs in each jaw of adults, without enamel, the dentine intersected by numerous tubules (hence the ordinal name). Four toes on forefeet, 5 behind, each bearing a strong, stout, hooflike "nail." Skull elongate and tubular; tongue vermiform; mandible not so reduced as in other "anteaters." Legs powerful, skeleton generally primitive but tibia and fibula fused. One species; total length to 200 cm (79 in.).

Order Lagomorpha (lagomorphs). Upper Paleocene to Recent of Eurasia, North America, Africa, northern South America. Dentition superficially rodent-like, but with 2 pairs of upper incisors (a small pair situated behind the larger functional pair); canine and anterior premolars lost, forming prominent diastema (gap); cheek teeth high-crowned, rootless, 6 above and 5 below (more than in any rodent), and not occluding directly, so that mastication is a lateral movement; maxillary bone fenestrate (*i.e.*, with holes). Fibula articulating with calcaneum; infraorbital canal small. Scrotum anterior to penis. Total length, 15 cm (6 in.; pikas) to 70 cm (28 in.; hares). Examples: hares, rabbits, pikas.

Order Rodentia (rodents). Lower Eocene to Recent, worldwide except Antarctica and some oceanic islands. Single pair of rootless incisors above and below; canines and anterior premolars lost forming prominent diastema; never more than 2 premolars above and 1 below; cheek teeth frequently rootless, hypsodont; cusp patterns variable, cheek teeth ranging from pegs to massive teeth with highly complex occlusal patterns. Powerful jaw musculature developed in connection with grinding dentition; angular process of mandible well developed for muscle attachment. Postorbital bar absent; infraorbital foramen variously modified, in advanced groups passing portions of masseter muscle; postcranial skeleton relatively unspecialized. Forelimbs flexible, digits clawed, little reduction in number on forefeet, some specialization of hindfeet and reduction of digits associated with saltation; stance plantigrade or semiplantigrade. Total length of living species, 10 cm (4 in.; pygmy mice, some pocket mice) to 125 cm (49 in.; capybara). Examples: mice and rats, squirrels, porcupines, beaver.

†*Order Creodonta* (creodonts). Fossil only; Middle Cretaceous to upper Oligocene of North America, Eurasia, Africa. Primitive carnivorous mammals; canines well developed; cheek teeth specialized for shearing, but carnassial pair farther back in tooth row than in Carnivora. Auditory bullae unossified. Distal phalanges channelled for insertion of claws; no fusion of scaphoid, lunar, and centrale bones of wrist. More than 50 genera, in about 5 families. Size to that of bear or larger.

Order Carnivora (carnivores). Middle Paleocene to Recent on all continents (probably accompanied early man to Australia). Teeth variously modified, sectorial to bunodont, but the carnassial (shearing) pair of teeth always consisting of the last upper premolar over the first lower molar; dentition simplified in advanced pinnipeds (seals), nearly homodont, with no carnassials. Auditory bullae ossified in advanced forms. Terminal phalanges not channelled for insertion of claws; scaphoid, lunar, centrale bones of wrist fused. Carnivorous, omnivorous, herbivorous. Total head and body length 15 to 370 cm (6–146 in.). Examples: cats, dogs, bears, weasels, civets, hyenas, raccoons, seals, sea lions, walruses.

†*Order Archaeoceti* (archaeocetes, Zeuglodonts). Fossil only; lower Eocene to middle Miocene of North Africa, Europe, North America. Primitive whale-like marine mammals, dentition not highly modified; anterior teeth simple, peglike; cheek teeth sectorial as in early carnivores; teeth not exceeding basic placental number (44); nostrils placed well back on elongate snout; skull not "telescoped" as in modern whales; hind limbs vestigial by upper Eocene; body elongate, some serpentine; total length to more than 2,000 cm (65 ft).

Order Odontoceti (toothed whales). Upper Eocene to Recent in all oceans and occasionally in freshwater. Dentition generally simple, undifferentiated (homodont), exceeding the primitive placental number (to 300 teeth in some porpoises) or secondarily reduced (to a single tooth in *Monodon*, the narwhal); nostrils located on top of head, above and between eyes, forming single "blowhole"; skull asymmetrical, markedly "telescoped," with reduction of frontals and parietals and concomitant extensive posterior expansion of maxillaries and premaxillaries. Relative to Archaeoceti, body shorter, stouter, total length 150 cm (59 in.; small porpoises) to about 1,900 cm (62 ft; male sperm whales). Examples: dolphins, porpoises, narwhal, sperm whales.

Order Mysticeti (baleen or "whalebone" whales). Upper Eocene to Recent in all oceans. Teeth absent; "whalebone"

(baleen) present, consisting of ridged "curtains" of keratin extending from roof of mouth, used to strain krill (zooplankton); mouth cavernous, esophagus constricted to maximum diameter of about 20 cm; mandibular symphysis lacking; skull symmetrical, markedly "telescoped," with anterior expansion of occipitals; nostrils moved posteriorly, "blowhole" double, anterior to orbits (eyes). Total length from 600 cm (236 in.; pygmy right whale) to more than 3,000 cm (98 ft; blue or sulfur-bottom whale); the blue whale, weighing to 136,000 kilograms (150 tons), is the largest of all known vertebrates, living or extinct. Examples: rorquals (including blue whale), gray whales, bowheads, right whales.

†*Order Condylarthra* (condylarths). Fossil only; Upper Cretaceous to lower Oligocene of Eurasia, North and South America. Probably the basal stock from which most later ungulate groups evolved; directly descended from primitive placentals, and characteristics diverse. Tendency toward rounded molar cusps of equal height (bunodonty), formation of diastema, reduction of canine teeth, assumption of digitigrade stance and "experimentation" with hooves (ungules). About 90 genera and 7 families; some species at least as large as modern bears.

†*Order Pyrotheria* (pyrotheres). Fossil only; Eocene and Oligocene of South America. Incisors enlarged, tusklike; diastema prominent; cheek teeth low-crowned and with transverse crests; nostrils on top of heavy skull, suggesting development of fleshy proboscis. About 4 genera known; as in many other ungulate groups, trend to large size (to that of elephants); length of skull to 60 cm (24 in.).

†*Order Xenungulata* (xenungulates). Known only from a few fragmentary fossils from upper Paleocene of South America. Generally poorly known, but molars with tendency toward converging transverse crests as in Dinocerata of Northern Hemisphere; postcranial remains scanty but build obviously massive. Size to that of a modern rhinoceros.

†*Order Pantodonta* (pantodonts). Fossil only; Paleocene to Oligocene of North America and Asia. A varied group of herbivores, prominent in the Paleocene. Cheek teeth low-crowned, with transverse crests variously arranged, canines small or greatly developed. About a dozen genera; total length approximately 100 to 450 cm (39–177 in.).

†*Order Dinocerata* (uintatheres). Fossil only; Paleocene and Eocene of North America and Eocene of Asia. Skull low, later representatives with bony projections on nasal, maxillary, and frontal bones. Upper incisors absent, lower incisors reduced, saber-like upper canines (evidently only in males); upper molars unique, with 2 transverse crests forming a "V." Limbs graviportal (massive); about 8 genera; size to that of large rhinoceros.

Order Proboscidea (elephants and allies). Upper Eocene to Recent in Eurasia, Africa, North and South America. Earliest known representatives characterized by low-crown molars with two transverse lobes, diastema resulting from loss of outer lower incisor, canine and first premolar, and formation of upper and lower tusks by modification of incisors. Nasal opening toward top of skull, indicating that a fleshy proboscis ("trunk") was present, but perhaps of modest size; early proboscideans relatively long-bodied and short-limbed, superficially tapir-like. Middle and later Cenozoic times saw wide radiation from basal stock; about 2 dozen fossil genera known. Two living genera (*Elephas*, the Indian elephant; *Loxodonta*, the African elephant) with high and truncate skull; massive tusks; short mandible; peculiar cheek tooth arrangement in which a single, high-crowned, lophodont tooth is exposed in each half of the jaw at any one time, and replacement occurs continually from behind; long proboscis; graviportal limbs and girdles; length of head and body to 640 cm (252 in.).

Order Sirenia ("sea cows," dugongs, and manatees). Middle Eocene to Recent of Africa, Asia, North and South America. In living species: hind limbs lacking, pelvis vestigial; front limbs modified into flippers; tail flattened, expanded; body fusiform, neck indistinct, anterior teeth reduced, replaced by horny plates used in cropping estuarine vegetation; molars of living dugongs bunodont, premolars lost; cheek teeth of manatees bilophodont, increased to 20 or more in each jaw, with 6 in each jaw at a given time and worn teeth replaced from behind as in modern proboscideans. Total length to 450 cm (177 in.).

†*Order Desmostylia* (desmostylians). Known only as fossils from lower Miocene to lower Pliocene in shallow marine deposits of North Pacific Ocean. Skull massive, low, elongate; anterior teeth modified as tusks, 1 pair (canines) above, 2 pairs (canines and incisors) below; cheek teeth peculiar, cusps forming cylinders of enamel and replaced in a continual forward movement as in elephants and manatees; limbs and girdles massive; appearance perhaps superficially hippopotamus-like; total length to about 250 cm (98 in.).

†*Order Embrithopoda* (embrithopods). Known only from fos-

sils from lower Oligocene of North Africa. Limbs 5-toed, progression graviportal; superficially rhinoceros-like, large paired horns on nasal bones and smaller horns on frontals (*i.e.*, behind the larger pair). Incisors not enlarged as in other subungulates; cheek teeth high-crowned; total length to about 330 cm (130 in.).

Order Hyracoidea (hyracoids). Lower Miocene to Recent of Africa and southwestern Asia. Most generalized of subungulates, although postorbital bar present and permanent dentition reduced; last upper incisors and canine lost, forming diastema; remaining incisors rootless; molars high-crowned. Superficially rabbit- or rodent-like; herbivorous, living on rock slopes or semiarid. Total length of fossil species to about 200 cm (79 in.); living members 30 to 65 cm (12–26 in.). Living examples: hyrax species.

†*Order Notoungulata* (notoungulates). Known only as fossils; upper Paleocene to Pleistocene of South America and Paleocene and Eocene of North America and Asia. Varied assemblage, united by dental peculiarities; accessory cusps in central valley of upper molars; lower molars with crescentic ridges, the entoconid lying within the arc of the larger posterior crescent; usually no diastema or reduction in dentition; cheek teeth tending to become rootless, high-crowned, lophodont. Tendency for reduction of toes from 5 to 3; total length approximately 30 to 300 cm (12–118 in.).

†*Order Astrapotheria* (astrapotheres). Known only as fossils; upper Paleocene to upper Miocene of South America. Upper incisors lost, canine enlarged, rootless; lower incisors and canine enlarged; diastema prominent; premolars reduced; molars greatly enlarged. Forelimbs stout, feet digitigrade, hind-limbs relatively weakly developed, semiplantigrade. Skull truncate; position and character of nasal opening suggests presence of fleshy proboscis; absence of upper incisors suggests that heavy lower incisors may have opposed a modified upper lip to allow cropping of vegetation. Total length to about 280 cm (110 in.).

†*Order Litopterna* (litopterns). Known only as fossils; upper Paleocene to Pleistocene of South America. Hoofed herbivores with tendency toward reduced toes; from 5 to 3 or even 1. Dentition usually complete (or with some loss of incisors); no diastema; cheek teeth low-crowned, selenodont; tendency for molarization of premolars. Size to that of small horse.

Order Perissodactyla (perissodactyls). Lower Eocene to Recent; North America, Eurasia, Africa. Incisors variously modified, canines lost or specialized, diastema nearly always present, with loss of first premolars above and then below; molars primitively bunodont, becoming progressively more lophodont, hypsodont; molarization of premolars common at later stages. Hornlike projections on skull not uncommon. Trends toward lengthening of limbs, with limitation of motion to forward and backward. Feet mesaxonic (*i.e.*, with axis through toe III), reduction typically to 3 digits and later to a single digit; digits mostly hooved, but those of chalicotheres clawed. Range in body size from that of small dog (Eocene horses) to that of largest of known land mammals, *Baluchitherium*, an Oligocene rhinocerotid of Asia, nearly 600 cm (236 in.) in height at the shoulder. Examples: titanotheres, chalicotheres, tapirs, rhinoceroses, horses.

Order Artiodactyla (artiodactyls). Lower Eocene to Recent; worldwide except Australia. Dentition complete, bunodont (as in pigs) or variously reduced; incisors often reduced and upper incisors lost, vegetation cropped with lower incisors opposing gums of premaxilla; canine and first premolars frequently lost, forming prominent diastema; remaining premolars simple, not completely molariform; cheek teeth primitively low-crowned, bunodont, becoming hypsodont, selenodont (cusps expanding into crescents) in more advanced forms. Unguligrade, foot paraxonic (*i.e.*, the axis between digits III and IV); digit I lost first; digits II and V frequently rudimentary (dew claws), but more prominent in pigs and hippopotamuses; in advanced types metapodials III and IV fused to form cannon bone, adding extra functional joint to limb; astragalus distinctive, the articular surface a double pulley allowing efficient flexion and extension of hindlimb but precluding lateral movement. Parietal bones of skull reduced, frontal bones expanded, frequently bearing antlers or horns. Height at shoulder, 25 cm (10 in.) in small antelope to 370 cm (146 in.) giraffe. Examples: pigs, cattle, sheep, goats, antelope, pronghorn, giraffe.

Critical appraisal. Three subclasses of Mammalia generally have been recognized for many years: Prototheria, the monotremes; Allotheria, the extinct multituberculates; and Theria, including marsupials and placentals as well as the extinct symmetrodonts and pantotheres. However, a number of Mesozoic groups could not be placed satisfactorily in any of these subclasses as defined. Evidence on

the formation of the braincase in nontherian mammals has suggested that the primitive taxa may in fact constitute a more coherent group than was once supposed. All nontherian mammals are thus placed in a single subclass, Prototheria; infra-ordinal status is proposed for monotremes, triconodonts and symmetrodonts, and multituberculates.

Monotremes and multituberculates. Monotremes and multituberculates represent highly modified offshoots of the earliest mammalian stock. The fossil record of undoubted monotremes is limited to the Pleistocene. Surely the group is an ancient one and, although living representatives are highly specialized, represents the grade of evolution of the reptilian-mammalian transition. Their primitive organization is underscored in the relatively unspecialized brain, retention of the oviparous habit and a cloaca, incomplete homeothermy, and reptile-like frontal and pterygoid bones.

The multituberculates also were derived from the earliest mammalian stock. They were adapted early to a herbivorous diet and developed gnawing incisors and elongate, grinding cheek teeth. The fossil record of the multituberculates, from Upper Jurassic to lower Eocene, is longer than that for any other mammalian order. The extinction of multituberculates in the Early Cenozoic may have been due to competition from placental herbivores, in particular the early rodents.

Marsupials. The split between marsupials and placentals has yet to be identified with any certainty. *Pappotherium* of the Lower Cretaceous of Texas was thought for a time to represent a common ancestor of the therian infra-classes, but more recently it has been discovered that recognizable didelphoid marsupials were contemporaneous with the pappotheriids.

The earliest records of marsupials are from North America but the great adaptive radiation of this group occurred in South America and Australia. Mostly free of competition from placentals on Cenozoic island continents, marsupials converged in many cases with the specializations of placentals elsewhere. Australia has been isolated from other major land masses throughout the Cenozoic and South America was an island from Paleocene to Pliocene. Invasions by placentals in later Cenozoic times reduced the once flourishing South America marsupial fauna to the remnant extant today. The South American and Australian radiations were independent, although both assemblages had a common source in generalized didelphoid (opossum-like) groups and evolved in parallel fashion in a number of instances. In general, higher categories in the classification of marsupials are not coordinate with those of the Eutheria. Some families of marsupials, for example, encompass a range of diversity comparable to that in eutherian orders.

Insectivores and colugos. Modern insectivores are divisible into several subgroups. The lipotyphlans include modern shrews, moles, and hedgehogs. The molars of hedgehogs are quadritubercular; those of shrews and moles are complicated by the development of a W-shaped external ridge or ectoloph. Lipotyphlans maintain the primitive insectivorous habits of their Mesozoic forebears. Zalambdodonts are represented by the peculiar selenodonts of the Greater Antilles, the tenrecs of Madagascar and West Africa, and the chrysochlorids—Cape golden moles—of southern Africa. The zalambdodonts may or may not represent a natural group. Their early history is obscure. All have V-shaped upper molars.

Tree shrews (tupaids) of Southeast Asia are the sole survivors of a third insectivore group, sometimes called prototherians, which includes a great variety of Late Cretaceous and Early Cenozoic mammals of uncertain relationships. Living tupaids have frequently been classed as Primates, for they have remarkably advanced brains and versatile limbs. In such features as dentition, however, they remain quite primitive.

The macroselids, or elephant shrews, are here treated as a separate order. They have frequently been grouped with tupaids as a suborder of the Insectivora, or as a distinct Order Menotyphla. Colugos, the so-called flying lemurs, here maintained as a distinct order, Dermoptera, likewise have been considered by some authorities as insectivores.

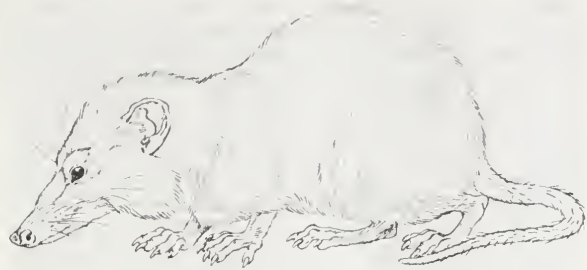


Figure 5: Reconstruction of *Melanodon*, a shrew-sized insectivorous pantothere from the Late Jurassic of North America.

Drawing by R. Keane

Tillodonts and taeniodonts. The relationships of tillodonts and taeniodonts are obscure. Although apparently quite independent, both groups were evidently derived from primitive placental stocks, and both evolved rootless, rodent-like incisors. Neither group was ancestral to later, more successful mammalian types.

Bats. Bats are an ancient group; except for extensive skeletal and sensory modifications associated with flight they share many primitive characteristics with their insectivore ancestors. Bats first appear in the fossil record in the Eocene as highly adapted fliers. Were well-preserved postcranial (body and limb) remains of the earliest Cenozoic insectivores better known, some might well be found to represent intermediate stages in the evolution of flight; to date no such intermediate stages have been found. Many have thought that the Dermoptera might represent a grade through which the lineage of the bats must have passed in evolving truly volant habits; surely dermopterans and bats are not closely related, save in their retention of primitive characters of early placentals. In recent years research has brought to light important differences between the two suborders of bats, the Megachiroptera and the Microchiroptera, and some workers have suggested that the two groups be treated as separate orders, maintaining that the similarities between them are due to parallelism and convergence in adapting to the stringent demands of flight.

Primates. The Primates form a relatively coherent group, although the distinction between them and their insectivore ancestors is a difficult one and some groups, such as the tree shrews, can be placed in either order with nearly equal justification. Disagreement on classification within the order centres on the validity of the suborder Prosimii as a natural unit that includes such divergent kinds as tarsiers, lemurs, lorises, and galagos, and on the relationship between the Old World and New World monkeys. The two monkey groups may have evolved independently from prosimian ancestors; considerable convergence and parallelism no doubt are involved in the similarities between the two as both have radiated widely in response to available arboreal habitats.

Edentates and pangolins. The edentates have undergone their greatest radiation in South America. They are a highly aberrant group, of unknown ancestry, and apparently were isolated early in South America, as were the archaic South American ungulates. Edentates have commonly been subdivided into Loricata (armadillos and extinct glyptodonts), Pilosa (tree sloths and extinct ground sloths) and Vermilingua (anteaters). Each of these subgroups has had a long and independent history, but they share the unique character of xenarthrous vertebrae. The palaeonodons of the Eocene and Oligocene of North America have been considered as nonxenarthrous edentates (a view presented in the section on *Edentata*); they are here allocated tentatively to the order Pholidota. Little is known of the relationships of pholidotes (pangolins), once considered relatives of the Edentata; the groups are now generally conceded to be convergent.

Aardvarks. Aardvarks (Tubulidentata) also were thought at one time to be related to true edentates, but are now considered an independent lineage. Although they are known from Miocene deposits, no direct evidence exists of relationships with other groups. In some respects the skeleton resembles that of certain early condylarths, and

it seems probable that tubulidentates descended from an early ungulate stock.

Lagomorphs. The Lagomorpha long were included (as a suborder Duplicidentata) in the Rodentia on the basis of dental similarities. Although ordinal rank for the Lagomorpha has been firmly established, some recent classifications still include rodents and lagomorphs together in a cohort Glires, again based on dental similarities, which almost certainly represent convergence. Various other mammalian lineages of diverse affinities (multituberculates, taeniodonts, tillodonts, and certain primates, for example) have "experimented" with gnawing incisors. The earliest known fossil lagomorphs clearly represent the group and yield no clues as to origin. On the basis of serology and some dental characters, a relationship to ungulates through the condylarths has been suggested. More than likely, however, those characters held in common are primitive and lagomorphs arose directly from a primitive insectivore stock.

Rodents. In terms of numbers of individuals and numbers and diversity of species, the rodents are the most successful of living mammals and have evidently been so throughout much of the middle and later Cenozoic. The earliest known rodents (paramyids), of Paleocene age, were already highly specialized mammals. Fossil evidence indicating relationships with more primitive groups is lacking; rodents may have been derived directly from a basal eutherian stock.

The higher classification of rodents has been in a state of flux in recent years. Three suborders (Sciuromorpha, Myomorpha, and Hystricomorpha) were recognized, based on the relationship of the zygomaseteric musculature to the infraorbital foramen. The sciuromorpha, or Protrogomorpha, included the earliest known rodents and also modern squirrels, beavers, and sewellels or "mountain beavers." In these taxa the masseter muscle originates from the anterior portion of the zygomatic arch or from the side of the skull in front of the orbit. The infraorbital canal is small, carrying nerves and blood vessels as in primitive mammals. In the Myomorpha, the rats and mice and their allies, the deep portion of the masseter muscle passes through the enlarged infraorbital canal. The so-called hystricomorphs, including Old World porcupines and a great variety of South American rodents, have a greatly enlarged and round infraorbital foramen, rather than a slitlike one like that of myomorphs.

Recent authors have doubted the validity of this arrangement and the utility of the infraorbital canal and zygomaseteric structure as a diagnostic feature. Many would admit the general validity of the Sciuromorpha and Myomorpha although some groups traditionally included would be re-allocated. The "Hystricomorpha," on the other hand, is almost unanimously considered to be polyphyletic, and is now commonly subdivided into Caviomorpha for the complex South American "hystricomorph" radiation, and Hystricomorpha for Old World porcupines and their allies. The tendency has generally been to de-emphasize the use of subordinal classification and instead to bring together related families into a number of superfamilies. Perfection of the higher classification of rodents is difficult because fossil evidence is meager for some groups, particularly those of small or modest size, and because the successful radiation of the Rodentia has obviously led to perplexing convergence in numerous lines.

Creodonts. Creodonts were derived from the basic insectivore radiation as early carnivorous mammals, largely in response to the ecological opportunity presented by the evolution of herbivorous types.

Carnivores. True carnivores evidently are not descendants of the creodonts, but direct derivatives of Paleocene insectivores. The diversification of fissioned (terrestrial) carnivores from the basal stock (family Miacididae) occurred in the upper Eocene and lower Miocene, early dividing into arctoid or canoid (canids, ursids, procyonids, mustelids) and aeluroid or feloid (felids, hyaenids, viverrids) lines. The Pinnipedia (seals, sea lions, walruscs) are here considered a suborder of the Carnivora. Appearing in the fossil record in the Miocene, the modern pinniped families already were distinct. The previous history of the group

is obscure, but pinnipeds apparently were derived from arctoid fissipeds, although some have thought them to have been derived from creodonts. Some authorities have treated pinnipeds as a distinct order, in which case they represent the most recently evolved mammalian order.

Whales. The whales and their allies have traditionally been treated as a single order, Cetacea, a reasonable arrangement on the basis of certain criteria. It is obvious, however, that cetaceans diverged early into archaeocetes, odontocetes, and mysticetes, if, indeed, the groups had a common ancestry. The earliest known fossils are perfectly good whales and clearly are assignable to one of the three groups; thus there is no concrete evidence linking the three nor is there evidence as to their ancestry. The morphological characters that whales share can be interpreted as due to parallelism resulting from stringent selection in similar habitats, or to the retention of primitive placental characteristics. Until evidence is available to establish the monophyly of cetaceans it seems preferable to treat them as representing three separate orders.

Ungulates. The term "ungulate" denotes a broad, loose association of orders including most large herbivorous mammals. The name of the group implies that its members are hoofed, but some of the ungulates (chalicotheres, for example) had clawed digits and others (such as sirenians) show considerable reduction of digits and limbs. The ungulates may be of common ancestry or may be a product of convergence from several sources among primitive placental stocks. Certain supposed early ungulates (condylarths) and early carnivorous mammals (creodonts) have many characters in common; their supposed descendants frequently have been united in a common cohort, Ferungulata.

The principal characteristics shared by ungulates are those of the dentition and limbs, related to modification of the primitive placental into an efficient herbivore.

There has been a tendency toward elongation of individual molar cusps to form crescentric ridges (selenodonty) or to form ridges between adjacent cusps (lophodonty). A concomitant trend has been an increase in the height of molars from low-crowned (brachyodont) to high-crowned (hypsodont), which has allowed the efficient utilization of abrasive forage such as grass. The extent of the grinding surface has been increased by increasing the length of the molars or by "molarization" of premolars.

The radiation of large herbivores was accompanied closely by a radiation of large carnivores. The adaptive response of ungulates to predation has been varied: modified dentition (tusks), skeletal elements (horns, antlers), and behaviour have been turned to defensive purposes. Extremely large size or swift, cursorial locomotion were at a premium. Defensive gregariousness is common in extant species.

The general trend in mammalian locomotion has been plantigrade to digitigrade to unguligrade, frequently with reduction of digits. One mode of digit reduction has been mesaxonic, the axis of the foot passing through digit III (as in horses); other groups have paraxonic limbs, with the axis lying between digits III and IV. Elongation of metapodials has occurred commonly, forming a third functional segment of the limb, allowing fast cursorial gaits.

Exceptions to these trends in limb modifications occur in lineages with early trends toward extremely large size. In such groups, columnar, graviportal limbs have developed with proximal segments (femur, humerus) emphasized and metapodials de-emphasized; all digits usually are retained.

The ungulate radiation has been complex. Various subdivisions have been proposed. The condylarths may or may not be basal to the entire radiation, but they evidently gave rise to the perissodactyls.

Perissodactyls. This is a compact order, consisting of five subdivisions: horses, chalicotheres, titanotheres, rhinoceroses, and tapiroids. The history of the perissodactyls is perhaps better known than that of any other mammalian group. *Hyracotherium* ("Eohippus"), usually considered the base of the horse lineage, was probably near the source of the entire order.

Artiodactyls. The Artiodactyla first appear in the early Eocene. Dental characters point to an origin for the group among hypsodontid condylarths but the earliest known artiodactyls already had the unique and specialized "double-pulley" astragalus. Knowledge of transitional stages between condylarthran and artiodactyl tarsi would be of considerable interest and importance.

Archaic ungulates. A number of kinds of archaic ungulates, which generally evolved in the Western Hemisphere, have been lumped in the past in a single order, Amblypoda. Each of these groups—pantodonts, uinatheres (Dinocerata), pyrotheres, xenungulates—is here accorded ordinal rank. The relationships of these groups—if any—are obscure and their origin is unknown.

Liopterns. The Liopterna of South America were early derivatives of the condylarths, and some workers would include both in a single order. One family (Prototheriidae) closely paralleled the contemporaneous radiation of perissodactyls on the northern continents, but the macrauchenids converged toward the camels.

Notoungulates. The order Notoungulata was an early and varied assemblage, apparently restricted to South America in the Eocene, perhaps by competition with early members of more advanced ungulate groups. The greatest development of the order was in the Oligocene, after which there was a gradual decline, hastened by the entry into South America of placental carnivores and modern artiodactyls. The origin and affinities of the notoungulates are unknown, but the order may well have descended from the condylarths.

Astrapotheres. Astrapotheres were once considered to have been notoungulates. It is now generally conceded that their peculiarities recommend recognition as a separate order. This group may have had its origin among primitive notoungulates, or it may have arisen independently.

Hyracoids. The products of the early Cenozoic radiation of ungulates in Africa are commonly grouped as "subungulates." The Hyracoidea is the most generalized

From A S Romer, *Vertebrate Paleontology*, 3rd ed. copyright 1966 the University of Chicago Press

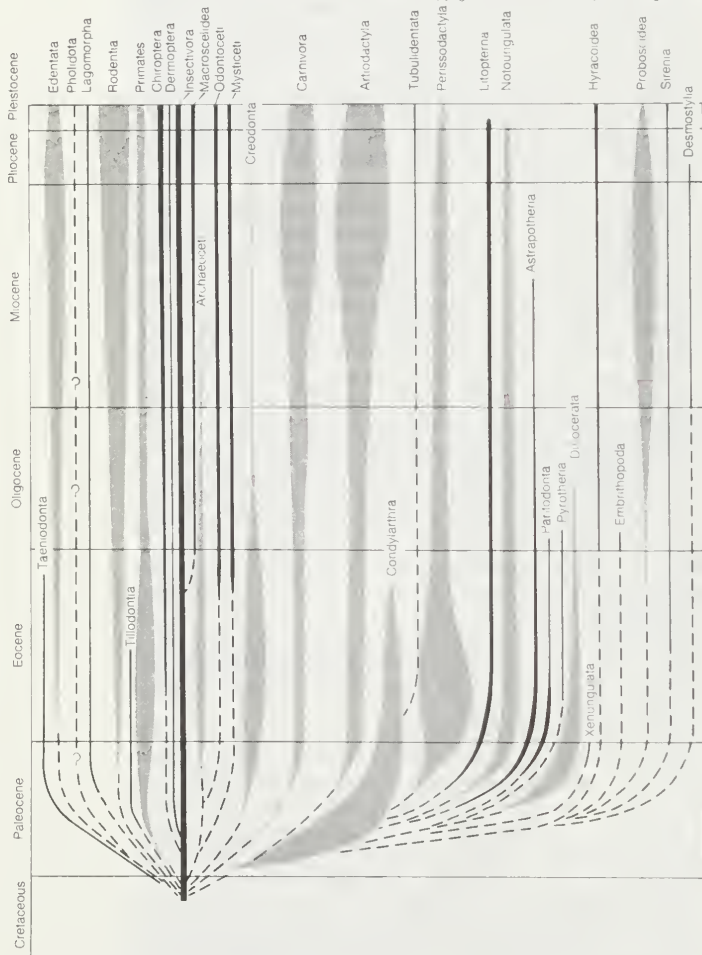


Figure 6: Dendrogram of the placental mammals.

of these orders but the true relationships of the hyracoids to other ungulates are unknown. In common with other subungulate orders, hyracoids were more diverse in the Middle Cenozoic than they are today.

Embrithopods. The extinct Order Embrithopoda contains a single known genus, *Arsinoitherium*. Relationship with hyracoids has been suggested but concrete evidence is lacking. Aside from logical association with other subungulate groups, little can be said of the position or affinities of the Embrithopoda.

Elephants and allies. The best known of the subungulates are the Proboscidea. The later record of the order is relatively good and indicates a remarkably successful and widespread radiation in the latter part of the Cenozoic—a radiation characterized by the recurrence of several trends in separate lines. Among these trends were an emphasis on the development of the upper tusk with loss of the lower tusk in most groups, increase in the size of the skull and shortening of the neck, increase in the length of the proboscis, increase in overall size, with perfection of graviportal limbs and girdles, lengthening of the lower

jaw (secondarily shortened in some groups), and reduction and specialization of cheek teeth.

Sirenians. The origin of the Sirenia is obscure. The earliest known sirenians were completely aquatic although the pelvic girdle was not so much reduced as in later forms. Dental and cranial peculiarities suggest that proboscideans and sirenians may have arisen from a common swamp-dwelling or amphibious ancestor. Concrete evidence is generally lacking regarding this and other relationships among subungulates due to a paucity of fossiliferous continental sediments of early Cenozoic age in Africa.

Desmostylians. Because of dental peculiarities and their association with shallow-water marine sediments, fossil remains of the order Desmostylia were long thought to represent aberrant sirenians; recent discoveries of massive limb bones of these mammals indicate, however, that they had quite different adaptations. The nature of the cheek teeth and their mode of replacement suggests relationships with the largely African subungulates, although desmostylians are known only from fossil deposits along the shores of the North Pacific. (J.K.J./D.M.A.)

MAJOR MAMMAL ORDERS

The following sections consist of a review of living mammals, beginning with the most primitive—monotremes and marsupials—and continuing with the major orders of placental mammals. Groups are identified at the order level of biological classification and arranged in accordance with the *Annotated classification* above.

Extinct groups and species (designated in the annotated classification by a dagger [†]) are treated throughout this article under the subheadings *Evolution and paleontology*; for additional information about extinct mammals, see *GEOCHRONOLOGY: Fossil record*, which also treats the subject of mammalian evolution. (Ed.)

Monotremata (platypus, echidnas [spiny anteaters])

The Monotremata (monotremes) are a distinctive order of primitive mammals, the only surviving members of the subclass Prototheria.

Monotremes lay eggs and are reptilelike in many other ways but possess such essentially mammalian characters as mammary glands, hair, a large brain, and a complete diaphragm.

GENERAL FEATURES

Only two monotreme types, the platypus (*Ornithorhynchus anatinus*, family Ornithorhynchidae) and the echidnas, or spiny anteaters (*Tachyglossus aculeatus* and *Zaglossus bruijnii*, family Tachyglossidae), are known. The earliest monotreme fossils known come from the Australian Pleistocene (2,000,000 or more years ago), and they are essentially the same as the living forms, which range in body length from about 30 to 80 centimetres (12 to 32 inches) and weigh from about one to 10 kilograms (two to 22 pounds).

The platypus, found in eastern Australia from north Queensland to Tasmania, is one of the most remarkable of all mammals. It has a ducklike "bill," webbed feet, and a flattened beaver-like tail. When the first stuffed specimens reached England around the end of the 18th century, they were thought to be fakes, made by sticking together bits of different animals.

The echidnas also rank high among the world's most interesting mammals and are readily kept in captivity for physiological studies. Superficially, they are quite unlike the platypus. They are recognized by their sharp-pointed spines and tubelike noses. *Tachyglossus* is found throughout Australia, including Tasmania, and part of New Guinea. *Zaglossus* has an extensive distribution in New Guinea.

Neither the platypus nor the echidnas are readily seen by the casual visitor, but they may be fairly common in some areas. Platypuses are not kept in zoos outside Australia.

The monotremes do not appear to be carriers of diseases harmful to man. Both monotreme types are now wholly protected by law in Australia.

NATURAL HISTORY

Ecology. The platypus is amphibious, with a habitat ranging from alpine waters at an altitude of 1,500 metres (about 4,900 feet) in southeastern Australia to the sluggish subtropical Queensland coastal streams. It digs winding burrows with side branches in the banks of the streams and lakes that it inhabits. These burrows are usually five to 10 metres (15 to 30 feet) long but may be as long as 30 metres (100 feet). The entrance is normally above the water level, however during floods it may be submerged.

The platypus has few natural enemies. It was formerly hunted extensively for its valuable velvet-like fur but for many years has been completely protected by law, and as a result it has increased in numbers. It appears to be relatively free of competition from other species in the water, where it feeds mainly on insect larvae, small crayfish, tadpoles, and other aquatic animals. Some authorities have suggested that the introduced rabbit (*Oryctolagus cuniculus*), which burrows in the banks of streams where the platypus once dwelt in quiet seclusion, is responsible

Enemies of the platypus

Drawing by A.G. Lyne

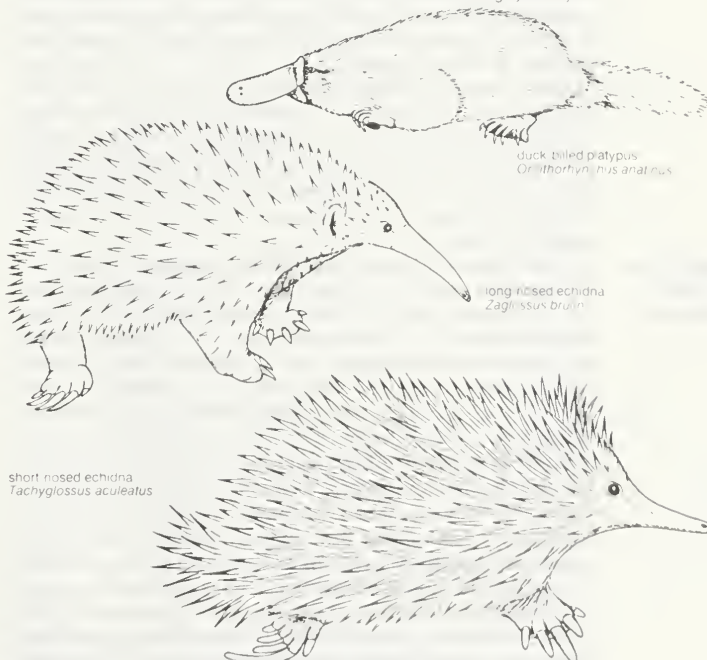


Figure 7: The three extant monotreme species.

for the disappearance of the platypus from many settled regions.

The echidnas are terrestrial animals, living in forests, scrubland, and deserts. Although usually encountered only by accident, *Tachyglossus* is one of the most widely distributed mammals in Australia. In New Guinea the habitats of *Zaglossus* include the humid forests that are almost continually blanketed by cloud cover. The animal has a known altitudinal range in these forests from about 1,100 to about 2,900 metres (3,600 to 9,500 feet).

Tachyglossus lives largely on ants and termites. The one element common in all the habitats occupied by this echidna is the presence of ants; termites are present in only some of the localities occupied.

Unfortunately, little is known of many aspects of the ecology of the monotremes, such as population density and social behaviour.

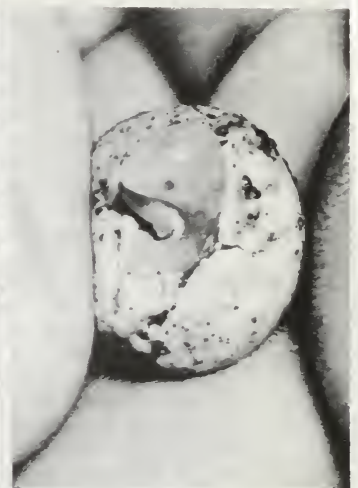
Life cycle and reproduction. The platypus breeds once annually, from July to August in the northern part of its range and to October in the southern part. Mating is believed to occur in the water, and it has been suggested that the spurs on the hind legs of the male may be used to hold the female. No spur grip was noted in an observation of copulation recorded by Australian zoologist D.H. Fleay in the course of an extensive study. An unusual courtship precedes mating, in which, among other manoeuvres, the male grabs the female's tail and the pair swim in circles.

In the breeding season, the female platypus digs a special burrow that twists and turns and may be very long. She then builds a nest of grass, leaves, and other plant material in a chamber at the end of the breeding burrow. This special burrow is blocked with earth in one or more places by the female during the period when she is laying and incubating the eggs. Only the left ovary is functional, and the eggs are always found in the left uterus. One to three, usually two, leathery-shelled eggs are laid. When there are two or three eggs they adhere together, but the sticky fluid that covers the eggs when they are laid soon dries and loses its stickiness. The eggs vary from 16 to 18 millimetres (about 0.60 to 0.70 inches) in length and from 14 to 15 millimetres in breadth.

The platypus has been bred in captivity on one occasion only, by Fleay in 1943. His observations suggest that the period from mating to egg laying is probably 12 to 14 days and the incubation period an additional 10 to 12 days. The young are not more than about 18 millimetres long at hatching and do not leave the breeding burrow until they are about 17 weeks old. At this age they are fully haired and about 34 centimetres (13 inches) long. Because the female platypus has no teats, the milk is apparently sucked from the skin as it exudes from the ducts of the mammary glands. The life-span of the platypus is said to be 10 to 15 or more years; at least one captive animal has lived more than 10 years.

More research has been done on the echidna *Tachyglossus* than on the platypus, mainly because it is less difficult to maintain in captivity. In southeastern Australia, echidnas breed once a year from late June to late September. The period between copulation and egg laying is at least 16 days and may be as long as 28 days. Both ovaries are functional.

In the breeding season the female echidna develops an abdominal pouch. Egg laying has not been observed, but there is circumstantial evidence that the egg is laid directly into the pouch by the female's curving her body during this process. One egg (rarely two or three) is laid at a time; it is creamy yellow and the shell has a rubbery texture. The eggs of echidnas are slightly smaller than those of the platypus, being from 14 to 17 millimetres (0.55–0.67 inch) long and from 13 to 15 millimetres wide. Within the pouch the egg is not stuck to the hairs but is quite loose and from time to time is located in different positions, apparently retained by the apposition of the lips of the pouch. The incubation period of the egg in the pouch is 10 to 10.5 days, and the newly hatched echidna is about 15 millimetres (0.6 inch) long and weighs about 0.38 gram (about 0.01 ounce). After hatching, the young lives in the pouch and imbibes milk by sucking two special areas of the skin—the mammary areolae—where the ducts



(Left) Echidna (*Tachyglossus aculeatus*) held head downward to show egg in pouch. (Right) Echidna hatching.
M.E. Griffiths

from the mammary glands open to the exterior on either side of the midline of the pouch. As in the platypus, no teats are present.

Authorities disagree on whether the young echidna takes milk by sucking or by licking. After detailed observations, an Australian zoologist, Mervyn Griffiths, stated that advanced pouch young take milk by sucking and not by licking it off the hairs in the mammary region, but he considered it possible that newly hatched and early stages of pouch young get their milk by licking it up as it is ejected onto the surface of the skin.

Information on how long the young stay in the pouch and what happens to them when they are cast out is meagre. One of the echidnas studied by Griffiths increased in body weight from 240 grams (about 0.5 pound) to 850 grams (about 1.9 pounds) in 43 days, a mean daily increase in weight of about 14 grams per day. The young was cast out of the pouch when it reached a body weight of 400 grams (about 0.9 pound). Thereafter, the mother's pouch regressed quickly, but she continued to suckle the young for some time. At the time when the young was dropped from the pouch, its eyes were open and the sharp spines projected from the skin. In wild echidnas it has been reported that the mother deposits the young in a burrow when its spines emerge and it is too large to be carried in the pouch.

Tachyglossus has been bred in captivity on only three occasions, once in Berlin and twice in Basel in Switzerland. In captivity *Tachyglossus* has lived up to 50 years, and *Zaglossus* up to about 30 years.

Behaviour and locomotion. Both the platypus and echidnas exhibit extreme specializations for feeding and locomotion.

The platypus emerges from its burrow for feeding usually in the early morning or late afternoon, but is sometimes in the water at other times of the day. When it is burrowing or walking on land the web of the front feet, which extends beyond the claws, is folded under the palms. The muzzle and the claws of the front feet are used in burrowing. The platypus swims gracefully and expertly, using the front feet. The hind feet, in conjunction with the flattened tail, are used mainly for stabilization. The tail also assists the platypus in diving. The animal often swims along the surface with only the upper part of the muzzle and a small part of the head and body above the water. When submerged, the eyes and ears are closed by a fold of skin. The sensitive muzzle is an important organ for guiding the platypus while it swims blind, often close to the bottom. It may stir up the mud or move stones to locate its food. It surfaces often to breathe and to chew the prey, which it stores in cheek pouches while it is submerged. It can rest near the bottom for a few minutes by wedging itself under stones or logs. In captivity, the platypus has an enormous appetite, devouring each day the equivalent of half or more of its own body weight.

Egg laying
in the
echidna

Swimming
by the
platypus

Echidnas are generally solitary. They appear to be more uncommon than they are because of the inaccessible nature of the places they choose to hide in and because of their remarkable ability to dig vertically into the ground and cover themselves with dirt when disturbed.

Observations on *Tachyglossus* in various parts of Australia indicate that its activity is closely related to the ambient temperature. In north Queensland, they are active primarily at night; but in the more temperate parts of eastern Australia, they are active both day and night. In desert regions they retire into crevices in rock outcrops during the day, where the humidity is higher and they are protected from intense sunlight.

Although echidnas are well adapted for digging, they do not excavate burrows, and they dig primarily to obtain food and escape enemies. They also retreat into hollow logs or under rocks or roots and wedge themselves into crevices and even into slight depressions in such a way that it is practically impossible to dislodge them. They can withdraw their appendages and erect their spines, as does a hedgehog. Echidnas walk with the legs fully extended so that the undersurface of the body is relatively high off the ground and with the hind toes directed outward and backward. They can run and climb well. *Tachyglossus* is extremely powerful and can tear apart rotten logs in search of termites or dig down to ant nests in rocky terrain.

FORM AND FUNCTION

General characteristics. One of the most striking features of the monotremes is that they lay eggs from which the young are subsequently hatched. These eggs are like those of reptiles in that they have large yolks and rubbery shells. Other interesting reptilian features are seen in the structure of the urogenital organs and in the skeleton. In both sexes of the monotremes, as in all reptiles, the intestine, bladder, and reproductive organs all open into a common chamber, the cloaca, with only one external opening. The penis of monotremes is used solely for delivering sperm and not also for urination as in most other mammals. In the skeleton, the shoulder girdle not only retains well-developed coracoid bones, which are reduced to vestiges in other mammals, but also an interclavicle. The pelvic girdle carries epipubic bones, which, as in some reptiles and marsupials, act as supporting structures for the ventral (abdominal) body wall.

The skull of monotremes has a smooth, rounded cranial portion, terminating in a long rostrum, or snout. Both monotreme types are toothless as adults, and the echidnas are toothless throughout life. The platypus has teeth when young, but they are shed before the animal becomes adult and are replaced by horny pads. The testes are abdominal and there are no teats and no vibrissae (whiskers). The limbs are modified for digging or swimming. The posture of the platypus resembles that of reptiles, especially lizards. In the echidnas, however, the limbs support the body well off the ground, even when the animals are stationary or walking slowly.

An interesting feature of the digestive system of monotremes is that the whole inner surface of the stomach is lined with a cornified (horny) epithelium. There are no glands of any kind. Particularly interesting is the lack of those that produce hydrochloric acid and peptic enzymes, which together initiate the digestion of protein in other mammals. Some other insectivorous mammals also have cornified epithelium covering most of the stomach lining, but there is a small glandular area. It has been suggested that the breakdown of food in the stomach of *Tachyglossus* is assisted by the grinding action of ingested dirt.

Platypus. The platypus has a flattened body covered with dense, short, fine hairs and coarse guard hairs, usually dark brown on the back and much paler on the belly. There is practically no neck. The tail is broad and flat like that of a beaver. It is covered with coarse hair, which is usually worn off on the under surface. The adult male is 50 to 60 centimetres (20 to 24 inches) in total length and the adult female is 40 to 50 centimetres (16 to 20 inches). An average-sized male weighs nearly two kilograms (4.4 pounds). One of the most conspicuous features of the platypus is its elongated and flattened muzzle, which bears

a superficial resemblance to a duck's bill. This muzzle is covered with a soft, rubbery, naked, and sensitive skin, which is a blue-gray colour on its upper surface and pale pinkish or mottled below. Dorsal (toward the back) and ventral flaps, or shields, of naked skin extend backward from the muzzle over the adjacent hair. The nostrils are located on the upper surface of the muzzle near its anterior end. The eyes are small and there are no external ears. Both the eyes and ears open into a facial furrow that is closed when the platypus is below the surface of the water. There are cheek pouches in which food is stored until it can be chewed. The limbs are short and the feet webbed. The webbing on the hind feet reaches to the base of the claws, and on the front feet beyond the claws, making extra large paddles for swimming.

The male platypus has a sharp, movable, horny poison spur on the inner side of each hind limb near the heel. Each spur is about 15 millimetres (0.6 inch) long and connected by a duct to a venom-secreting gland, called the crural, or poison, gland, situated on the dorsal aspect of the upper part of the hind limb. The function of the venom apparatus is unknown. An Australian zoologist, J.H. Calaby, stated that "If, as is generally believed, the crural gland is actively secreting only during and near the breeding season it seems very likely that the chief use of the venom apparatus would be combat between males for territory or females." The poison is not fatal to man but causes intense pain.

Echidnas. Superficially, the echidnas are quite unlike the platypus. They have a rounded body covered on the back and sides with stout, sharp-pointed spines and some hairs. The spines of *Tachyglossus* may measure up to six centimetres (2.5 inches) in length, those of *Zaglossus* up to about three centimetres. The undersurface is covered with coarse hair and is usually without spines. The colour of the spines and hair ranges from light brown to almost black, depending on species and individual variation. Adult males are larger than adult females. There is a long, naked, and sensitive muzzle, a very short tail, and no neck. In some echidnas a definite external ear is present, but in others it is inconspicuous. The nostrils open near the tip of the muzzle. The mouth is small, only sufficient for the protrusion of the long, sticky tongue. The limbs are short and powerful and have strong claws for digging. In *Tachyglossus* one of the claws of the hind foot is particularly elongated and is used for scratching the skin between the spines.

The echidnas resemble the platypus in that there is a venom apparatus on each hind leg, but the spur and crural gland are smaller. In *Tachyglossus*, the spur (about five to 10 millimetres [0.2–0.4 inch] long), which is situated on the inside of the ankle, is present in all males and some females. No observations on the use of the venom apparatus by the echidnas have been published.

Temperature regulation and hibernation. Monotremes are better able to control their body temperature than are reptiles. Except during hibernation, or torpidity of some species, the eutherian, or placental, mammals maintain a high and virtually constant body temperature usually between 36° and 40° C (97° and 104° F). The body temperature of most reptiles varies with that of the surroundings. The body temperature of the monotremes, however, does not quite reach that of the eutherian mammals, being about 31° to 32° C (88° to 90° F). The monotremes are usually described as physiologically primitive in regard to temperature regulation. It has been shown, however, that *Tachyglossus* is an excellent temperature regulator in some situations but less capable in others. It is probably no poorer in this respect than many eutherian mammals (e.g., anteaters, armadillos, and sloths). In cold surroundings the echidna maintains a well-regulated body temperature. Studies of temperature regulation and torpor in *Tachyglossus* have shown that this echidna appears to be a true mammalian hibernator, with the ability to arouse from torpor without the aid of heat from external sources. This ability is lost after repeated periods of torpor in the laboratory at 5° C (41° F). The echidna resembles other hibernators in that torpor is not continuous but is interrupted by frequent arousals. Continuous periods of torpor

Venom apparatus of the platypus

Similarities to reptiles

lasting 9.5 days, with intervening periods of wakefulness ranging from 30 hours to 11 days, have been recorded for the echidna.

The platypus has periods of hibernation or semihibernation, at least in southeastern Australia. In captive animals, Fleay found that the hibernation periods were short and irregular during the cooler months, the longest single period being 6.5 days.

Laboratory studies have shown that the platypus is slightly superior to *Tachyglossus* in temperature regulation. The platypus has sweat glands on the muzzle and over the body surface, whereas *Tachyglossus* has them only in the pouch area. The heat tolerance of *Zaglossus*, which has many well-developed sweat glands distributed over the body, has not been studied. In its natural environment it is unlikely that the platypus will be exposed to high heat loads. The highly insulating coat of the platypus would protect it against heat loss in water and allow it to exploit cold environments.

EVOLUTION, PALEONTOLOGY, AND CLASSIFICATION

The monotremes are without living relatives outside Australia and New Guinea, and, apart from the presence of fossil monotremes in the Australian deposits of the relatively modern Pleistocene Epoch (the last 2,500,000 years or so), little is known of their ancestry. The fossils so far discovered differ from the existing forms only at the species level. It is interesting to note, however, that large echidnas of the genus *Zaglossus*, which now live only in New Guinea, were common and widespread in Australia during the Pleistocene.

Although the existent monotremes are highly specialized in many of their features, their possession of such typically mammalian characters as mammary glands, hair, and a large brain carries the suggestion that such characters are primitive features of early mammals.

Most authorities regard the order Monotremata as the sole group of the subclass Prototheria and believe that it originated from a line of mammal-like reptiles different from that which gave rise to the other mammals. It has been suggested, because of the similarities between marsupials and monotremes, that they should be grouped together in a subclass of the Mammalia, the Marsupionta, but the idea has not gained wide acceptance. It does not take into account the many similarities between marsupials and eutherian mammals. It has also been postulated that a group of early prototherian mammals, of which the monotremes are specialized derivatives, gave rise to the marsupials and the eutherian mammals, but this view is not accepted by most biologists.

The ducklike muzzle, or "bill," alone distinguishes the platypus from all other mammals. Although platypuses have been divided into several subspecies, based mainly on differences in size, all authorities agree that the family Ornithorhynchidae has only a single genus and species, *Ornithorhynchus anatinus*.

Two genera of echidnas are recognized in the family Tachyglossidae; in *Tachyglossus* (short-nosed echidnas) the elongated muzzle is straight and not greatly prolonged, and in *Zaglossus* (long-nosed echidnas) the elongated muzzle is greatly prolonged and curved downward.

Although Griffiths recognized *Tachyglossus aculeatus* as the only species of the genus, he distinguished six subspecies, some of which he considered of doubtful validity. Some authorities consider that the relatively hairy Tasmanian echidna is a distinct species (*T. setosus*). Some mammalogists who have studied *Zaglossus* in detail have concluded that only one variable species, *Z. bruijnii*, should be recognized. It can be seen that much more work on the biology of the monotremes needs to be done to clarify the place these extraordinary animals occupy in nature.

(A.G.Ly.)

Marsupialia (kangaroos, bandicoots, phalangers, opossums, koala, wombats)

The Marsupialia is an order (or superorder, depending on the authority) of mammals characterized by premature birth and continued development of the newborn while

attached to the nipples on the lower belly of the mother. The pouch, or marsupium, from which the group takes its name, is a flap of skin covering the nipples. Although prominent in many species, it is not a universal feature among marsupials; in some species, for example, the nipples are in a well-defined area but are fully exposed or are bounded by mere remnants of a pouch. The young remain firmly attached to the milk-giving teats for a period corresponding roughly to the latter part of development of the fetus in the womb of a eutherian, or placental, mammal.

The largest and most varied assortment of marsupials—more than 100 species—is found in Australia alone: kangaroos, wallabies, wombats, the koala, and a bewildering assemblage of smaller rodent-like forms. About 70 more species are distributed more widely, in Australia (including Tasmania), New Guinea, and a cluster of nearby islands. The wide array of Australian marsupials is reflected in the extensive popular vocabulary of names, many of which are derived from descriptive Aboriginal words. Only two families of marsupials—totalling more than 70 species—are found in the Americas, vestiges of a larger group that originated there as long ago as the Cretaceous Period (from 136,000,000 to 65,000,000 years ago). The family Didelphidae comprises about 65 species of South and Central American opossums, one of which ranges as far north as southern Canada. The family Caenolestidae consists of seven species of ratlike marsupials confined to South America.

GENERAL FEATURES

Scientific and economic importance. Marsupials are of interest to zoologists for several reasons: their current geographical distribution, their remarkable evolutionary expansion in Australasia, their similarity in many respects to placental mammals, and, most obviously, their reproductive adaptations and often bizarre structure.

Fossil evidence indicates clearly that the marsupials originated in the New World, and although the oldest fossils referable to marsupials are found in North America, it is believed that South America is equally likely as the origin of these animals. The current concentration of diverse types of marsupials in Australia and its nearby islands is thought to have occurred as a result of passage over presumed land connections with South America during a mild early geological period before the rise of the placental mammals. Later, toward the beginning of the Tertiary Period (starting 65,000,000 years ago), Australia was isolated from all other continental masses, and the marsupials were free to evolve unrestricted by competition from other groups of mammals. Elsewhere, however, the marsupials did badly; they faced strong competition from and were supplanted by the more advanced placental mammals.

Australian marsupials provided products—meat and hides particularly—for the Aboriginal people, whose primitive hunting methods posed no threat to the continued success of the animals as a group. With the appearance of modern methods of hunting and trapping, however, several species of the kangaroo family (Macropodidae) were soon rendered extinct, and many others were brought close to that fate. For a time during the first half of the 20th century, many of the larger marsupials were slaughtered in great numbers for their pelts, which were an item of export in the fur trade; for their hides, which were made into shoes; and for their flesh, which was processed into dog and eat food. The gray kangaroo (*Macropus giganteus*), for example, is a casualty of such wholesale killing. Fortunately, in Australia most marsupials have since that time been accorded protection under law, and efforts are proceeding to conserve many species that were brought to the brink of extermination. The Tasmanian wolf, tiger, or thylacine (*Thylacinus cynocephalus*), thought to have become extinct in the 1930s, may still survive. In 1961 an expedition claimed to have found thylacine tracks in a rain forest on the west coast of Tasmania. The Tasmanian devil (*Sarcophilus harrisi*), once dangerously close to extinction, is now thriving in Tasmania. In the Americas marsupials are also hunted locally for food and other products, but, since numerous and equally desirable placental mammals abound, hunting pressure and the resulting threat of species death are not so strongly felt.

Origin of the order

The question of origin and spread

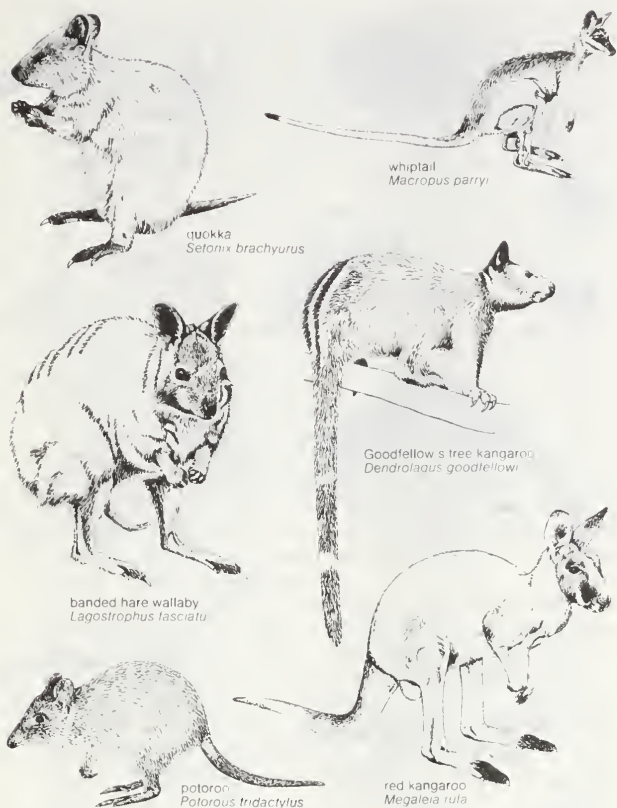


Figure 8: Representative species of the family Macropodidae. Drawing by A.G. Lyne

Not only man but also his introduced animals have played a part in the extinction of certain marsupials. Foxes, dogs, and cats continue to prey upon many species of marsupials, as also do the dingoes, canines introduced by the Aborigines long ago and since gone wild. Brush and forest fires as well as climatic change and the encroachment of civilization are also factors tending to reduce the numbers of marsupial species.

Size range and diversity. No extant marsupial even approaches the tremendous size of certain extinct forms, which apparently rivalled the mastodons in bulk. Among the largest living marsupial is *Macropus giganteus*, some individuals of which reach about two metres (six feet) in height, three metres (10 feet) from muzzle to tail tip, and weigh up to 90 kilograms (about 200 pounds); red kangaroos (*Megaleia rufa*) grow to about the same size. The smallest is the planigale (*Planigale ingrami*), a species of marsupial mice that is barely 12 centimetres (4.7 inches) in total length. The vast majority of marsupials lie in the range from the size of a squirrel to that of a medium-sized dog (Figures 8-10).

Structural and behavioral parallels with placental mammals are in some cases quite striking. Such resemblances are examples of convergent evolution, a tendency for organisms to adapt in similar ways to similar habitats. Thus, there are marsupials that look remarkably like moles, shrews, squirrels, and mice. Others, less in structure than in habits, are the ecological counterparts of cats, small bears, and rabbits. Even the larger grazing marsupials, which resemble no placental mammal at all, can be thought of as filling the same ecological role as deer and antelope found elsewhere.

Distribution and abundance. There are about 175 species of marsupials native to the Australasian area, primarily in Australia (including Tasmania), New Guinea, Timor, and Celebes. Through man's agency, however, marsupials have been introduced to nearby islands of Australia and especially to New Zealand. The more than 70 native American species occur in South and Central America, with the range extended into North America by the common, or Virginian, opossum (*Didelphis marsupialis*). A number of other species extend into the southern portion of North America.

In terms of numbers of species the didelphids, or New World opossums, rank highest with 65. In sheer abundance, there are more vegetarian than carnivorous marsupials and perhaps more pouched mice than any other kind of marsupial.

The brush-tailed possums (*Trichosurus vulpecula*) are examples of marsupials that have readily adapted to changing conditions brought about by man. As recently as 1932, more than 1,000,000 pelts of this species were sold under such names as Adelaide chinchilla. Protected for part of the year, the brush tails have now become plentiful over the vast interior of Australia and even pestiferous in urban centres. Their adaptability to different locales is attributed to their tolerance for a variety of food, including household refuse.

Koalas (*Phascolarctos cinereus*), reduced to a few thousand in the 1930s as a result of disease and logging and trapping for fur, have recovered their numbers sufficiently to be reintroduced into forested regions.

Although the red kangaroos (*Megaleia rufa*) are being killed in large numbers for commercial products and by farmers and ranchers whose grazing lands they sometimes visit during periods of drought, they are numerous in many parts of inland Australia. The total range of this

Drawing by A.G. Lyne



Figure 9: Representative species of the families Phalangeridae, Petauridae, Vombatidae, Phascolarctidae, Peramelidae, and Tarsipidae.

Parallels with placental mammals

species is about 2,000,000 square miles. Changes in vegetation resulting from domestic stock have improved the habitat for it in many parts of its range, and the species has increased in abundance.

NATURAL HISTORY

Life cycle. The life cycle of marsupials exhibits the peculiarities of a mammalian group advanced over the egg-laying monotremes but primitive to the placental mammals. The uterine cycle of the female marsupial has no secretory phase, and the uterine wall is not specialized for the implantation of the embryos. The period of intrauterine development in marsupials ranges from about eight days in the native cat *Dasyurus viverrinus* to 40 days in the red-necked wallaby (*Wallabia rufogrisea*). The young, born in a vulnerable embryonic condition, make their own way to the shelter, warmth, and nourishment of the pouch; in pouchless marsupials the young simply cling to the teats. Those fortunate enough to survive this arduous journey may succeed in attaching themselves to the mother's nipples, which then swell and become firmly fastened—almost physically fused—to the mouth tissues of the young. In this condition the young continue their development for weeks or months, after which they are weaned and begin to look after themselves. Frequently the partially developed young outnumber the available teats, and the excess individuals perish. In *Didelphis marsupialis*, litters of as many as 25 young have been reported, but the pouch usually contains only 13 teats.

The primitive reproductive cycle of *Didelphis marsupialis*

Didelphis marsupialis exemplifies one of the more primitive cycles and serves to illustrate the sequence of events in greater detail. It may have two litters yearly throughout most of its range (possibly only one in the northern reaches). The female experiences estrus (heat) and becomes receptive to the male about every 28 days during the breeding season. Heat lasts one or two days, during which an average of 22 ova are shed. Usually only 10 young are born in a litter.

During the first five or six days embryos form very slowly, but thereafter the rate of development speeds up. The gestation period is short—about 13 days. Just before birth, the mother thoroughly cleans the pouch and her belly fur, licking a smooth path from the birth canal to the pouch. When expelled from the birth canal, the newborn are no larger than honeybees. Blind and grublike, except for their well-developed, clawed forelimbs, they emerge and immediately begin grasping and swimming-like movements along the mother's dampened fur. The strong ones reach the pouch in less than a minute, an extraordinary feat for so small and incomplete a creature. Only about 60 percent of the newborn reach the pouch; these attach themselves immediately to the teats, at which they remain for four or five weeks (Figure 11). They then begin to leave the pouch for short intervals, scampering back to it whenever danger threatens. The young stay with the mother for 90 to 100 days—toward the last not in the pouch but clinging to her fur. Shortly afterward they are weaned and begin their independent lives.

Behaviour. The marsupials are notably less intelligent than placental mammals, a fact that is attributable in part to a simpler brain (see below *Form and function*). It is not surprising, therefore, to find a repertory of behavior that differs somewhat from that of the more advanced placentals. One peculiarity that may stem from this underdevelopment is restricted vocal ability. Although marsupials are not entirely silent, few of them emit loud sounds of excitement or distress; apparently, none utter grunts of contentment or even cries of hunger when young. What vocalizing they do is more limited and less variable than that of placentals.

There seems to be little detectable social organization among marsupials beyond the short-lived pair bonds during mating. The reproductive cycle in many marsupials is seasonally controlled, and estrus occurs in a predictable rhythm. In the quokka (*Setonix brachyurus*) and many other kangaroos and wallabies, however, the females may bear young at any time of the year. Only one young, called a joey, can be carried by the mother at a time. After suckling for about six months in the pouch, the

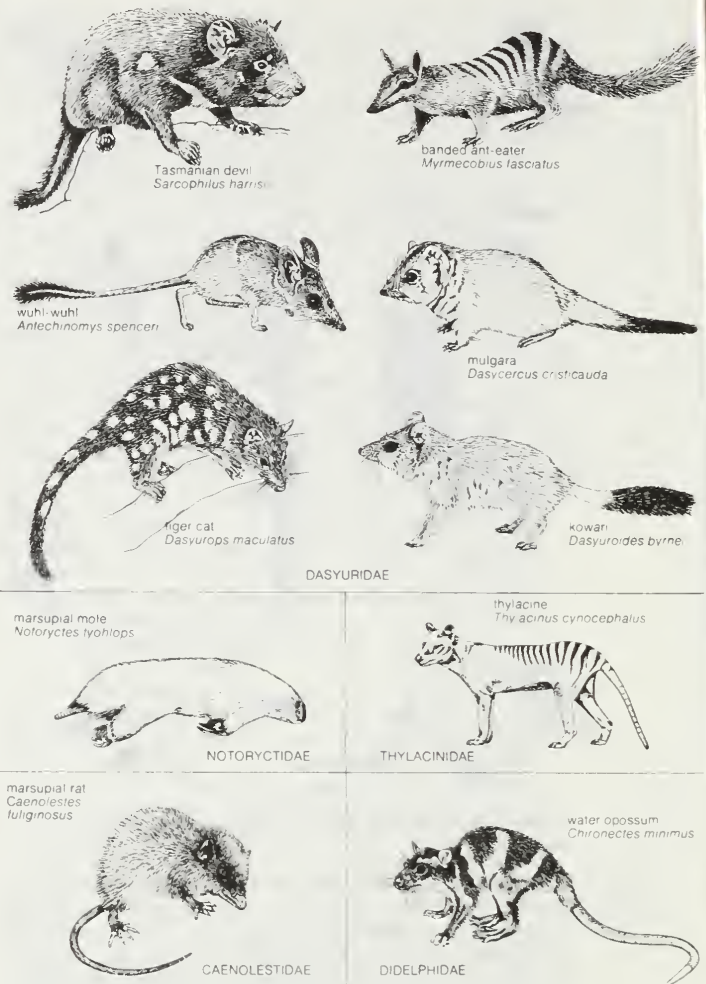


Figure 10: Representative species of the families Dasyuridae, Notoryctidae, Thylacinidae, Caenolestidae, and Didelphidae. Drawing by A.G. Lyne

joey begins to venture out but remains very close to the mother, plunging headfirst into the pouch when alarmed. The males have no interest in the young and will protect neither them nor the females against danger.

Many of the grazing marsupials, such as the kangaroos and wallabies, move in feeding groups called mobs, but these associations do not merit distinction as true herds since the individual members appear to move at liberty, without attention to any leaders or elders. One member can send the mob into a wild rout—individuals bounding off in all directions—by thumping its tail on the ground in a signal of alarm.

The larger marsupials, especially, embody nervous alertness and dull-witted consternation at the same time. A typical response of a red kangaroo when startled is to bound off in full flight—up to 65 kilometres per hour (about 30 miles per hour)—for a relatively short distance, stop short to reconnoitre the disturbance, and, in what appears to be sheer panic, double back toward its adversary and jump over it, sometimes clearing more than three metres (10 feet) in height and a distance of more than eight metres (25 feet) in a single leap.

Although well equipped for escape, a boomer, or large male kangaroo, can stand fast against many of its foes. With agile arms it can spar vigorously, behaviour learned in playful bouts as a joey. But, much more effectively, it can use the forepaws to grip its enemy, while it rocks back on its tail and then swiftly drops its huge clawed hindfeet, an action that has been known to disembowel dogs and men. Despite such fearsome ability, the kangaroos are usually gentle, as are most of the vegetarian marsupials. Whatever violence they commit is usually in response to a distinct threat to themselves.

Such marsupials as the possums (families Burramyidae,

Defensive behaviour of kangaroos

Petauridae, and Phalangeridae), wombats (Vombatidae), and koala (Phascolarctidae) usually sleep by day and move about toward dusk in search of food. Except for the agile possums—notably the gliders (*e.g.*, *Acrobates*, *Schoinobates*)—most of these animals are slow-moving and, consequently, easy prey for carnivorous animals. The bush-tailed possum, however, is not easy prey for carnivorous animals.

The most mild-mannered marsupial is perhaps the numbat, or banded anteater (*Myrmecobius fasciatus*), a small, slow forest prowler who searches out termites, its favourite food. When disturbed, the numbat hisses; if caught it neither struggles nor bites but simply vents a few grunts in protest. Equally inoffensive and similarly vulnerable to predation is the teddy-bear-like koala, which rummages about at night among eucalyptus branches munching endlessly on the only food it eats—eucalyptus leaves. Its name, an Aboriginal dialect word meaning “no drink,” refers to the koala’s curious and apparently lifelong abstinence from water. Koalas provide a eucalyptus “soup” for their young at the time of weaning; it is lapped up directly from the parent’s anus. The densely furred cuscuses (*Phalanger*) move among the trees with the same deliberation of the koalas, eating leaves as they go.

Bandicoots—whose apt name means “pig rat”—are especially fond of earthworms and insects and are only occasionally herbivorous. Although generally shy and retiring, bandicoots can at times display a surprising belligerence: with their sharp-clawed feet they can literally scratch their victims to death. Certain marsupial mice (*e.g.*, *Sminthopsis* species) are so hyperactive—like shrews—that to supply their high energy needs they must devour their own weight in food, chiefly insects, each day.

A curious and effective defensive behaviour has been developed to a high degree in *Didelphis marsupialis* of North America: when harassed it hisses and departs, but when attacked and not able to flee, it feigns death (plays possum) so well that the aggressor may lose interest and leave.

The carnivorous marsupials are equipped with other behavioural traits. They are the fierce hunters of flesh. Although as a rule not as fast as comparable placental carnivores, they are persistent in the hunt and swift to kill. The Tasmanian devil is a stout fighter, likened in its bloodlust to the American wolverine. Yellow-footed marsupial mice (*Antechinus flavipes*) cling so tenaciously to their prey that they can be caught like fish on bait. The thylacine, which has been known to rise up on its hindlegs kangaroo fashion when pressed, makes up in perseverance what it lacks in speed and cunning. In sharp, foxlike snaps, thylacines have been known to dispatch an antagonist as formidable as a hunting dog.

Ecology. The niches, or ecological roles, that marsupials fill are closely associated with structure. The burrowing species, such as the marsupial mole and the wombats, have powerful foreclaws with which they can tunnel into the ground for food and for shelter. Terrestrial forms, such as the kangaroos and wallabies, possess well-developed hindlimbs that serve both as formidable weapons and catapults by which they can bound over the plains. The greater glider (*Schoinobates volans*) and other Australian “flying” possums have a membrane along either flank, attached to the forelegs and hindlegs, that enables these arboreal animals to glide down from a high perch. A few marsupials spend most of their lives in trees; koalas, for example, are so thoroughly arboreal that they seem comfortable only when they are among the branches of trees. One marsupial is semi-aquatic, the water opossum—yapok or yapó (*Chironectes minimus*)—of Central and South America. Its thick, oily fur, partially webbed feet, and constrictable pouch opening suit it admirably for swimming and diving in search of food.

The diet of marsupials is as varied as the niches. Many live chiefly on insects and other small animals. Such are the marsupial mice and many of the smaller native cats (family Dasyuridae). The tiger cat and the Tasmanian devil feed largely on birds and small mammals. The numbat uses its remarkable wormlike tongue to lap up termites and ants. Many Australian possums, bandicoots, and American opossums have a mixed diet of plant matter and insects. Certain other marsupials are strictly vegetarian. The wombats feed on grasses, roots, and fungi, which they dig up with their clawed forepaws. The small honey possum (*Tarsipes spenserae*) is specialized to feed on the nectar of flowers.

The entire family of macropodids—the “big-footed” ones—are primarily grass eaters. They include the kangaroos and wallaroo, the wallabies and pademelons, and the rat kangaroos and musk kangaroo. The wallaroo, or euro (*Macropus robustus*), in particular, has a diet almost as restrictive as that of the koala: it can subsist in large part on spiky spinifex grass. Macropodids digest their vegetable food in the same manner as do sheep, cattle, and other ruminants, by relying on intestinal bacteria and protozoans to break down plant matter.

In their turn, marsupials are food for other animals. Foremost among the predators of didelphids and small dasyurids are owls and other birds of prey, snakes, and many carnivorous mammals, including dogs and cats. Some of the larger species of American opossums are eaten locally by man. The enemies of the larger marsupials, the macropodids, include the dingo, fox, eagle, pythons, goannas, and especially man.

Charles Philip Fox



Figure 11: Development of the common opossum. (Top left) In the pouch at three weeks; (bottom left) at seven weeks; (right) mother carrying 10-week-old young.

FORM AND FUNCTION

Marsupials share with other mammals the presence of hair and mammary glands. With minor differences the various systems of the body, such as the muscular and skeletal systems, are those of the placentals generally. The skull and the brain, however, differ considerably from those of placentals. Differences also exist in the dentition and in the arrangement of digits of the feet.

Bodily adaptations. In gross structure, many marsupials have hindlimbs that are larger and more powerful than the forelimbs; such an arrangement is most obvious in the large terrestrial grazing marsupials but is of wide occurrence in some degree throughout the Marsupialia. Arboreal and burrowing marsupials have forelimbs almost as well developed as the hindlimbs, understandable adaptations to their modes of living. Claws, often of considerable size, are invariably present in all marsupials and assist in climbing or digging. The tail, of tremendous size in kangaroos and wallabies, serves as an organ of stability and balance; in smaller and arboreal species it may also be prehensile (adapted for grasping), allowing the animal to hang by it.

Brain
peculiarities

Compared with placentals, marsupials differ markedly in both the structure and bulk of the brain. Most notably they lack a corpus callosum, that part of the placental brain that connects the two cerebral halves. In addition, the marsupial brain is smaller relative to skull size: a marsupial cat has about half as much brain tissue as a placental cat of similar skull size. Such insufficiency may account in part for the relative backwardness of marsupials when contrasted with placental mammals.

A feature peculiar to marsupials is the presence of a pair of bones associated with the pelvic girdle. Called the epipubic, or marsupial, bones, they were earlier thought to be a supportive element for the pouch.

The teeth in marsupials are numerous—usually 40 to 50, but as few as 22 in the honey possum and as many as 52 in the numbat. The relationship of milk teeth to adult teeth is not clear; some authorities claim that a single set of teeth lasts throughout life. Typically, there are seven cheek teeth on each side of the jaw, top and bottom, which are divided into three premolars and four molars; this contrasts with the typical placental condition in which there are four premolars and three molars. Also unlike the placental condition, the number of front, or incisor, teeth differs in upper and lower jaws.

In earlier studies of marsupials much was made of the number and type of incisors (Figure 12). Marsupials with more than three incisors in the upper jaw were said to be polyprotodont and of a presumed more primitive type since the teeth were more or less unspecialized. Marsupials with three or fewer upper incisors and a corresponding set of modifications of the primitive condition were said to be diprotodont. Subsequent investigations disclosed a gamut of intermediate dentitions, thereby discrediting the distinction.

Another equally unsatisfactory distinction was based on the curious but minor condition in which the second and third toes of the hindfeet are covered in a common sheath of skin. Species showing that condition were syndactylous, and the remainder were didactylous. Classifications were founded on such matters, and the groups were named appropriately Polyprotodontia and Diprotodontia in one scheme, and Syndactyla and Didactyla in another (see below *Classification*).

Reproductive adaptations. The most extraordinary anatomical features of the marsupials are the specializations associated with the reproductive system. The most striking of these is the protective pouch around the nipples of the female in many species. The pouch is well developed in kangaroos, wallabies, and wombats; in the tuan (*Phascogale tapoatafa*) and some dasyurids it is poorly developed, represented by lateral abdominal folds. In several South American didelphids, in the rat opossums (family Caenolestidae), and in the numbat, the pouch is lacking entirely. Even among species that have pouches, the degree of development depends upon the breeding condition of the animal, being most obvious when the female is lactating. The pouch in most forms opens forward, especially

Differences
in pouches

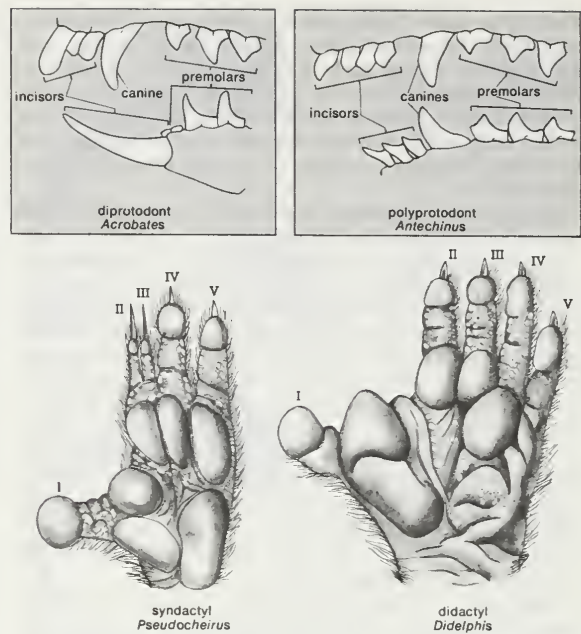


Figure 12: Dentition and fusion of the digits of the hindfoot of marsupials. (Top left) Diprotodont in the pygmy gliders and (top right) polyprotodont in the broad-footed pouched mice. (Bottom left) Digits II and III syndactyl, digit I opposable, in the ring-tailed opossums and (bottom right) digits II and III didactyl, digit I opposable, in the common opossum.

By courtesy of Owen J. Poe

in upright and arboreal species; in the water opossum and many quadrupedal species the pouch is directed backward. There may be as few as two nipples (as in the koala and wombats) or as many as 27 (as in *Monodelphis*, the short-tailed opossum). The number of teats does not indicate the number of young that are nursed, however; in some cases more young are born than can be suckled at one time, but in most cases fewer young than the number of nipples are brought to full term.

The female reproductive tract is double for most of its length, becoming joined at the posterior ends of the two lateral vaginae. A birth canal, or median vagina, forms at the time of birth; copulation is accomplished by the insertion of the male's penis into the urogenital sinus of the female. The end of the penis is forked in some species. In addition to the peculiar penis, the male also has a scrotum in which the testes hang in front of the penis rather than behind; the testes are abdominal in the pouched male (*Notoryctes*). Most marsupials specialize in brief gestation, as indicated earlier. The developing embryo obtains its food chiefly from its own yolk sac, an embryonic organ provided by the mother. In the majority of marsupials, embryonic gaseous exchange is also effected by means of the highly vascular yolk sac. Typically, there is no intimate mother-fetus connection by means of a true (chorioallantoic) placenta, as in placental mammals. The bandicoots are unique among living marsupials in having a rudimentary placenta, but it lacks the extensive ramifying projections (villi) that provide the close connection of mother and fetus tissue in the placentals.

EVOLUTION AND PALEONTOLOGY

Geologically, no marsupials have been discovered in Asia, and they are definitely not known from Africa. On the whole they are considered as primitive mammals, but for some 100,000,000 years they have shown a remarkable parallel and convergent evolution with placental mammals in habits, physiological processes, and structures.

Among the oldest marsupials are those represented by undisputed fossils from the Late Cretaceous strata of North America. These specimens are represented mostly by tiny isolated teeth, but parts of jaws have also been found. Until recently, all of the fragments were classified in the American opossum family, Didelphidae, but it has been observed that these early marsupials were already specialized in different directions. *Alphadon*, a genus of primarily

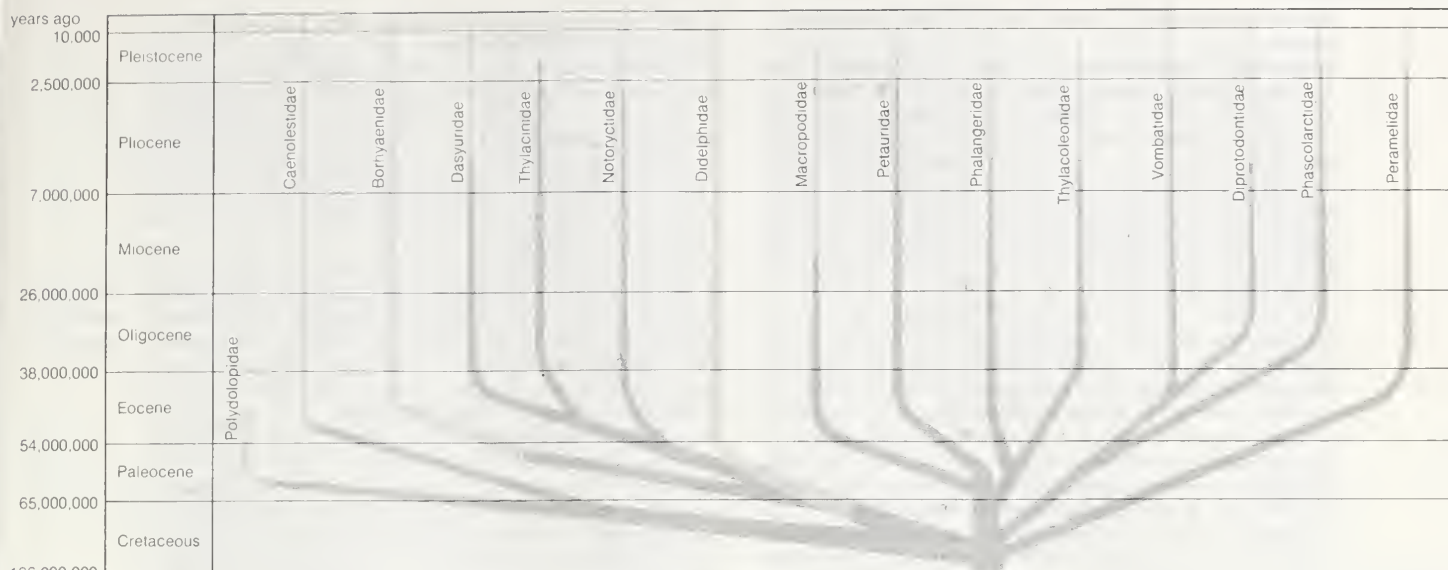


Figure 13: Dendrogram of the marsupials.

tiny animals, is most likely a didelphid. On the other hand the group of small to medium-sized creatures, represented by the genus *Pediomys*, seem rather clearly not didelphids and are now placed in their own family, the *Pediomyidae*. The widely known genus *Eodelphis* is almost certainly not didelphid. Among its more conspicuous differences from the *Didelphidae* are the single- (not double-) rooted, much reduced first lower premolar and the presence of three, not four, lower incisors, of which the middle incisor is greatly enlarged. It is likely that this marsupial is not ancestral to any of the known later genera. Another specimen, known as *Thlaeodon*, is the largest known Cretaceous mammal. Its bulbous premolars and large canines are indicative of an adaptation to a mixed (omnivorous) diet. It apparently is referable to another family. *Alphadon*, earlier thought to be nearer than any of the other Cretaceous genera to an ideal ancestral position for the marsupials as a whole, has been superseded by later finds. Fossil fragments of even earlier date (the Middle Cretaceous of Texas, about 100,000,000 years ago) have been tentatively referred to as *Holoclemensia* and assigned as a new genus in the family *Didelphidae*.

Evidence suggests that marsupials became widespread in the Tertiary Period, beginning 65,000,000 years ago, but they never succeeded as well on other continents as they did in South America and Australasia.

Primitive marsupial characteristics

American marsupials. *Didelphoids.* Primitive marsupial characteristics in the superfamily *Didelphoidea* have long been recognized by mammalogists. Some of the features in the dentition, skull, and body skeleton may be reasonably expected to be remarkably like those in the Mesozoic ancestors of all later marsupials. Characters in the molars indicate that didelphoids and the South American borhyaenoids of the Paleocene Epoch (beginning 65,000,000 years ago) and Eocene Epoch (beginning 54,000,000 years ago) are closely related, but the lack of enough fossils makes it impossible to trace the lineage of all marsupial superfamilies back to the known North American Cretaceous fossils (possibly the more direct ancestors of some of the other superfamilies existed at that time). Better knowledge of the skeletons and dentitions of fossil marsupials from the Cretaceous Period and even the Early Cenozoic Era (beginning 65,000,000 years ago) may eventually alter considerably the present interpretation of their superfamily relationships.

Borhyaenoids. The oldest known South American marsupials are of the later Paleocene Epoch (about 55,000,000 years old). At that time three superfamilies, *Didelphoidea*, *Borhyaenoidea*, and *Caenolestoidea*, were already in existence. Patterns of the known cheek teeth of the Paleocene and early Eocene borhyaenoids and didelphoids indicate that these marsupials may have descended from a common ancestor no earlier than the Late Cretaceous.

The fossil record of the family *Borhyaenidae* starts in

the late Paleocene and continues into the Pliocene (ending 2,500,000 years ago). No genera are known to have occurred outside South America. *Borhyaenids* were the flesh-eating mammals of the southern continent during nearly all the Tertiary Period. (True carnivores presumably did not reach that region until late in the Pliocene and Pleistocene.) The family name is derived from the genus *Borhyaena*, hyena-like specimens found in the early Miocene strata (about 26,000,000 years ago) of Argentina. These carnivorous marsupials had massive skulls with thick crushing teeth. Some genera—e.g., *Proborhyaena*—were shorter faced and more robust.

All borhyaenids were not hyenoid, however. There were the wolflike *Prothalycinus* and *Lycopsis* and the fox- to marten-sized *Cladosictis* and *Amphiproviverra*. One of the most specialized of all was the Pliocene sabre-toothed borhyaenid *Thylacosmilus*: it was about the size of a puma and even more specialized as a stabbing mammal than most sabre-toothed cats. In all, more than 20 borhyaenid genera have been described, revealing a long evolutionary history in South America from a Cretaceous didelphid-like mammal.

Necrolestoids. One fossil with a molelike adaptation, called *Necrolestes*, is known from the Miocene of Patagonia. It is the only known genus that can be referred to the family *Necrolestidae* (superfamily *Necrolestoidea*) and is probably basically related to the *Borhyaenoidea* or *Didelphoidea*.

Caenolestoids. The South American rat opossums (superfamily *Caenolestoidea*) are also undoubtedly descendants of a Cretaceous group, but the oldest known caenolestoids are already too specialized to reveal affinities with any of the known Mesozoic fossils. Any resemblances they have to Australian groups apparently result from convergent evolutionary trends in certain structures.

The early rat opossums

The family *Caenolestidae* has been recorded from early Eocene deposits, but the genera were more numerous later in the Oligocene and Miocene periods. These little marsupials were long known from fragmentary fossils, but in 1895 this supposedly extinct family was discovered as represented in the living Andean mammalian fauna of Ecuador and Bolivia by an animal resembling a small rat in size and appearance and named *Caenolestes*. The diprotodont incisor teeth are suggestive of those in some of the Australian phalangeroids, but the hindfeet show no trace of the phalangeroid syndactylism of the second and third toes. Other genera with living species are *Orolestes* in Peru and *Rhyncholestes* in Chile.

The *Polydolopidae* is an Early Cenozoic family of small specialized marsupials related to the *Caenolestidae*. Their fossil remains are abundant in the South American strata of that time. They had transversely compressed and laterally ridged premolars, with serrated crests, and rodent-like incisors. The polydolopids, probably in a large measure,

occupied the ecological niches that were later occupied by rodents.

Australasian marsupials. Although pre-Pleistocene fossils of Australian marsupials are relatively rare, representatives of the families Dasyuridae, Peramelidae, Phalangeridae, Thylacoleonidae, Macropodidae, and Diprotodontidae appear as fossils in strata as early as the late Oligocene Epoch (26,000,000 years ago). During the Miocene and Pliocene epochs that followed, fossils of the family Vombatidae were laid down. These findings suggest that the Australian families of marsupials originated in the Early Cenozoic.

Dasyuroids. The oldest dasyuroid lived during the Middle Tertiary of South Australia; about the size of the present-day *Dasyurus quoll*, it was probably near the stem from which the thylacine evolved. *Glaucodon ballaratenensis*, one of the most important fossil dasyurids, is possibly of Pliocene age; it bore a resemblance to both *Dasyurus* and a Pleistocene member of *Sarcophilus*. The thylacine (*Thylacinus cynocephalus*), which constitutes the family Thylacinidae, is the largest and most widely known of the dasyuroid marsupials. Although its distribution is usually recorded as being restricted to Tasmania, it was represented on the mainland of Australia less than 10,000 years ago and also lived on New Guinea. The last certainly known thylacine was captured in Tasmania in 1930; however, sightings have been recorded of thylacines and their tracks as late as 1961 along a wild stretch of the west coast of Tasmania.

The doubt-fully extinct thylacine

The family relationships of the thylacine are still not fully understood. (Its classification with the extinct Borhyaenidae of South America has almost conclusively been disproved.) Its characters, however, display a remarkable convergent evolution with the American borhyaenids. It seems reasonable to assume that the thylacine and the dasyures arose from a common ancestry. Many authorities place *Thylacinus* in the family Dasyuridae, but until its ancestry is adequately known it seems best to treat it as forming a separate family.

Without a fossil record it is not possible to conclude how closely *Notoryctes*, constituting the family Notoryctidae, is related to the basal stocks of the dasyuroids and to that of the bandicoots, but the absence of syndactyly in the hindfeet suggests they are closer to the dasyuroids. *Myrmecobius*, also lacking a fossil record, is grouped with the family Dasyuridae because of similarities.

Perameloids. Members of the superfamily Perameloidea display one of the most primitive characters seen in any Australian marsupials: the retention of an incisor formula of $\frac{3}{3}$ (polyprotodont) in most of the genera. They have syndactylism in the hindfoot: the proximal and median phalanges (bones) of the second and third digits are enclosed in the same sheath of skin.

Among the perameloids, both *Macrotis* and *Perameles* are represented in the late Pleistocene faunas of Australia. An extinct genus, *Isechnodon*, has been found in the Pliocene strata of South Australia. The most peculiar of all peramelids is the probably extinct (it has not been seen since about 1926) pig-footed bandicoot (*Chaeropus ecaudatus*). The common name is derived from the construction of the front foot, in which the second and third digits are of equal length and closely united. The nails, of equal size and length, give the appearance of cloven hooves.

Phalangeroids. Only a very hypothetical phylogeny, indicating possible ancestral relationships of the superfamily Phalangoidea to the other superfamilies, can be outlined until an adequate fossil record from the Late Mesozoic and Early Cenozoic has been compiled. The Phalangoidea were probably derived from an ancestral stock that also gave rise to the Perameloidea. Evidently, in this early stock (in contrast to the Dasyuroidea) natural selection was toward syndactylism. Subsequently, one syndactyl stock led to the ancestors of the rat kangaroos, kangaroos, diprotodontids, wombats, koalas, marsupial lions, and phalangers; the other stock gave rise to the bandicoots.

Of the phalangers, *Eudromicia* is the most primitive of all the living genera. It comprises in part the dormouse possums. The retention of the fourth molars and the double-rooted first and second upper premolars distinguish *Eu-*

dromicia from the other dormouse possums, *Cercartetus*, which they notably resemble.

The remains of a tiny marsupial called *Burrarnys* has been found in New South Wales, Australia. The shape of the skull and the pattern of the molars resemble those in *Cercartetus* (including *Eudromicia*), but the premolars are high and serrated. This genus is the basis for the family Burrarnyidae. Once thought to be extinct, several specimens of *Burrarnys* have been found since 1966 in northeastern Victoria.

One of the oldest fossil marsupials of Australasia is part of the skeleton of a brushtail-like possum, *Wynyardia bassiana*, found in late Eocene or Oligocene deposits near Wynyard, Tasmania. It constitutes the family Wynyardiidae. Fossils more closely related to later species occur in the Pleistocene of Australia; others much older occur in the Middle Tertiary of South Australia.

The oldest Australasian fossil

Among the most bizarre of all marsupials are the woolly possums, or cuscuses (*Phalanger*), from which the family and superfamily names are derived. *Wyulda*, the scaly-tailed possum, which is terrestrial and lives among rocks in northwestern Australia, is somewhat intermediate between the brush-tailed possums and the cuscuses. There has been much discussion about the relationships of the so-called marsupial lion (*Thylacoleo*). Most authorities agree that it belongs in a separate family, the Thylacoleonidae. Specimens of the cranium and lower jaws are well-known, and parts of the limb bones and body skeleton have been found. *Thylacoleo* was about the size of an African lion. Its teeth were even more specialized for meat shearing than those in the true cats; for example, the large posterior premolars were as long as 5.7 centimetres (2.24 inches). All the remaining teeth were greatly reduced, except the upper and lower median incisors, which apparently were utilized in capturing other animals. These great carnivorous marsupials became extinct in Australia less than 26,000 years ago.

Although the genera referred to as the rat kangaroo family (Potoroidae) are closely related to the kangaroo family (Macropodidae) and are classified by some as a subfamily of that group, they apparently have represented a distinct lineage since Eocene time. They may be distinguished from kangaroos by, among other features, several peculiarities of dentition and skull features. The female urogenital system is more specialized than in the Macropodidae. A Pleistocene genus, *Propleopus*, is known, and a species of *Bettongia* and an undescribed genus occur in the Middle Tertiary of South Australia.

The most widely known marsupials are the kangaroos and their relatives (family Macropodidae). The gigantic forms during the Pleistocene were *Procoptodon*, *Sthenurus*, *Protemnodon*, and two species of *Macropus*. *Prionotemnus* is from the Pliocene, but the oldest kangaroos come from rocks probably as old as Miocene.

The macropodids are adapted for jumping. The principal digit in the hindfoot is the fourth, the fifth being somewhat reduced; the syndactyl second and third are reduced to splinters but still bear tiny claws. *Procoptodon* and *Sthenurus* differ from the other kangaroos in having lost all but the fourth digit.

Vombatoids. The superfamily Vombatoidea includes the diprotodonts, the koala, and the wombats.

The extinct family Diprotodontidae included the largest of all marsupials, *Diprotodon optatum*. It was built something like a huge ground sloth and attained the size of a large rhinoceros. Special dentition included the large, chisel-shaped median incisors and molar teeth, somewhat like those in certain kangaroos or tapirs, composed of two cross crests. *Diprotodon* was probably abundant in Australia in the late Pleistocene and subrecent time. Its remains have been found in levels contemporaneous with man. At least two other genera, *Nototherium* and *Euowenia*, representing much smaller diprotodontids, also occur in Australian Pleistocene faunas. *Palorchestes*, also of the Pleistocene, previously considered as the largest of all kangaroos, is now known to be a diprotodontid.

The massive diprotodonts

This family, like the Macropodidae, had a long Cenozoic history in Australia. *Meniscophus* is known from the Pliocene, and remains of a rather primitive undescribed

genus and species are thought to be as old as the Oligocene Epoch. It appears that all known diprotodontids were herbivorous. The family may have been an early offshoot of the phalangeroid stock and perhaps was distantly related to kangaroos and wombats.

One of the most famous of all Australian marsupials is the koala, family Phascolarctidae. Only one species is known, although a rather distantly related genus, *Perikoala*, from the Tertiary of South Australia, has been described. The koala, with its short tail and rather stocky body, so resembles the wombat that the two are considered by some authorities to have arisen from a common ancestor, perhaps not very far back in Cenozoic time. Others believe, however, that the koala is much closer to the ringtails and greater gliders. The wombats (family Vombatidae), which resemble koalas in having no tails, are represented by two Pleistocene fossil genera: the giant wombat, *Phascolonus*, which was as large as a black bear, and *Ramsayia*. (Ed.)

CLASSIFICATION

Distinguishing taxonomic features. Two earlier classifications of marsupials were based on the number and arrangement of the front teeth and on the appression or separateness of certain toes. Marsupials having at least four upper incisors were named Polyprotodontia; those having three or (usually) fewer upper incisors and no lower canines, Diprotodontia. Marsupials having entirely separate toes on the hindfeet were called Didactyla; those having the second and third toes of the hindfeet enclosed in a common envelopment of skin, Syndactyla. These terms no longer have taxonomic validity, but their adjectival forms are still used to describe the dentition and the character of the hindfeet.

Additional aspects, such as the pattern of hair tracts, blood chemistry, details of reproductive anatomy and physiology, chromosomal morphology, sperm morphology, and behaviour, are currently being employed to structure a new taxonomic system representing more natural units that reflect what are presumed to be fundamental phylogenetic relationships.

Annotated classification. The time-honoured position of the Marsupialia as an order of mammals, generally unacceptable for a variety of reasons (see below *Critical appraisal*), is gradually giving way to a system first presented in 1964. As further modified, that arrangement raised the Marsupialia to the rank of a superorder and encompassed four orders, 10 superfamilies, and 23 families, as given below. Wholly extinct groups are preceded by a dagger (†).

SUPERORDER MARSUPIALIA (Metatherian, or pouched mammals)

Origin during the Cretaceous, expansion at beginning of the Tertiary, and decline with the rise of placentals thereafter. Currently established in Australia and nearby islands (introduced on New Zealand) and in the Americas. Young born after brief gestation in an embryonic condition; development completed while young are attached to teats, usually in a pouch, or marsupium. Skull elements differ from placentals; braincase small; brain relatively small and simple compared with placentals. Epipubic bones associated with the pelvic girdle. Fewer incisor teeth in lower than in upper jaws. Female reproductive tract doubled in large part (paired vaginae and uteri); testes usually in a scrotal sac and anterior to the forked penis (abdominal testes in *Notoryctes*). Comprises about 80 genera and about 240 species in 4 orders.

Order Marsupicarnivora

Comprises 3 extinct superfamilies and 2 extant superfamilies of primarily carnivorous marsupials.

†Superfamily Argyrolagoidea

Comprises a single family, Argyrolagidae, represented by *Argyrolagus* from the upper Pliocene to the lower Pleistocene in South America.

†Superfamily Borhyaenoidea

Comprises a single extinct family, Borhyaenidae, flesh-eating marsupials of South America during the late Paleocene and into the Pliocene Epoch (about 5,000,000 years ago). Many species had massive skulls and heavy crushing teeth. Named after hyaena-like specimens from Argentinian fossils. Examples: *Prothalcinus*, *Borhyaena*, *Lycopsis*, and the sabre-toothed marsupial *Thylacosmilus*.

Superfamily Dasyuroidea

Primarily terrestrial carnivores that resemble the didelphids in certain anatomical features. Incisors small and unspecialized (polyprotodont). Digits in hindfeet never joined (didactylous). Marsupium not present in all species; when present it is poorly developed and opens backward. Three families ranging from Australia to nearby islands.

Family Dasyuridae. A group widespread over Australasia and represented by marsupial mice and rats (e.g., *Sminthopsis*, *Antechinus*, *Antechinomys*, *Planigale*, *Murexia*, and *Phascogale*), native cats (*Dasyurus*, *Dasyurops*), Tasmanian devil (*Sarcophilus*), numbat (*Myrmecobius*). Includes the smallest known extant marsupial, *Planigale ingrami* (about 12 centimetres [4.7 inches] in total length). About 18 genera and 48 species.

Family Notoryctidae (marsupial moles). Mole-sized marsupials, of the deserts of central and western Australia, remarkably convergent with placental moles of other continents. Pouch opens backward, and there are epipubic bones. One genus, *Notoryctes*, and 2 species.

Family Thylacynidae (thylacines). German-shepherd-sized wolflike marsupials. Represented on the mainland of Australia less than 10,000 years ago and in New Guinea, the thylacine was last captured in Tasmania in 1930. A single species, *Thylacynus cynocephalus*, constitutes the family.

Superfamily Didelphoidea

Primarily terrestrial American marsupials. Earliest fossils from Cretaceous strata. One extant family, Didelphidae, and 2 fossil families.

Family Didelphidae (American opossums). Central and South America boasts many species of unusual adaptations, notably *Chironectes* (water opossum), *Philander* (4-eyed opossums), *Metachirus*, and *Lutreolina*. *Didelphis* ranges as far north as southern Canada. Includes 12 genera and about 65 species.

†**Family Pediomyidae.** One genus, *Pedionys*, from the Upper Cretaceous in North America.

†**Family Stagodontidae.**

†Superfamily Necrolestoidea

Comprises 1 family, Necrolestidae. Based upon molelike fossils from the Miocene (starting 26,000,000 years ago) of Patagonia. One genus, *Necrolestes*.

Order Paucituberculata

Comprises a single superfamily of marsupials.

Superfamily Caenolestoidea

South American terrestrial marsupials. Earliest fossils from the early Eocene Epoch (about 50,000,000 years ago). One extant family, Caenolestidae, and 2 extinct families.

Family Caenolestidae. Ratlike in size and appearance. Diprotodont incisors suggest relationships to Australian phalangeroids, but other evidence disputes such affinities. Includes *Caenolestes* (Ecuador and Bolivia), *Orolestes* (Peru) and *Rhyncholestes* (Chile). Comprises 3 genera and 7 species.

†**Family Groeberiidae.** One genus, *Groeberia*, from the late Oligocene.

†**Family Polydolopidae.** Known from several genera including *Polydolops* and *Amphidolops* from upper Paleocene to late Eocene in South America.

Order Peramelina

Comprises 1 extant superfamily with 1 family of primarily carnivorous marsupials. Specimens of modern species are common in Pleistocene cave deposits in various parts of Australia.

Superfamily Perameloidea (bandicoots)

Comprises 1 family, Peramelidae. Terrestrial Australian marsupials that resemble rodents, rat-sized to hare-sized. Most genera have primitive polyprotodont dentition: 5 incisors above and 3 below. Syndactylism in hindfeet. Pouch is directed backward. Eight genera and 22 species, including: *Peroryctes*, *Microperoryctes*, and *Rhynchomeles*, primitive genera restricted to New Guinea and adjacent islands; *Perameles* and *Isodon*, of wide distribution; *Macrotis* (bilbies) and *Chaeropus* (pig-footed bandicoot), restricted to Australia. Largest of the family (weight to 7 kg, or 15 lb.) is *Peroryctes broadbenti*.

Order Diprotodonta

Comprises 3 extant superfamilies of primarily herbivorous marsupials.

Superfamily Phalangeroidea

Australasian marsupials ranging from squirrel-sized arboreal species to large terrestrial bounders. Syndactylism of the second and third digits in the hindfeet. Four extant families; 2 extinct families.

Family Burramyidae. Primarily arboreal mouse- to squirrel-sized marsupials. Includes *Acrobates* (feathertail gliders), *Bur-*

ramys, and *Cercartetus* (pigmy possums). Comprises 3 genera and 6 species.

Family Macropodidae. Primarily terrestrial medium- to large-sized Australian marsupials. Adapted for jumping, with long hindlegs and long tail for balance. Forelimbs have sharp claws, and thumb digit in hindfoot is the fourth. Extinct giant forms occurred during the Pleistocene, but oldest forms can be traced to the Miocene. Extant macropodids include *Macropus* (gray kangaroos and wallaroos), *Megaleia* (red kangaroos), *Wallabia* (wallabies), *Thylogale* (pademelons), *Dendrolagus* (tree kangaroos), and *Setonix* (quokkas). There are about 19 genera and 47 species.

Family Petauridae. Terrestrial and arboreal marsupials. First and second digits of the forelimbs are opposable to the other digits. Molars are adapted for chewing leaves. Includes *Pseudocheirus* (ring-tailed possums), *Petaurus* (sugar gliders), *Schoinobates* (greater gliders), *Gymnobelideus* (Leadbeater's possum), and *Dactylopsila* (striped possums). About 5 genera and 25 species.

Family Phalangeridae. Squirrel-sized to cat-sized arboreal species. Includes *Trichosurus* (brush-tailed possum), *Phalanger* (cuscuses), and *Wyluda* (scaly-tailed possum). About 6 genera and 15 species.

†**Family Thylacoleonidae.** So-called Australian marsupial lions, which became extinct less than 26,000 years ago. Teeth remarkably specialized for meat shearing. One genus, *Thylacoleo*.

†**Family Wynyardiidae.** Brushtail-like possums, the oldest fossil marsupials of Australasia, from presumed late Eocene rocks of Tasmania. One species, *Wynyardia bassiana*.

Superfamily Tarsipedeoidea (honey possums)

Comprises a single family, Tarsipedeidae. Adapted for feeding on nectar of flowers. One species, *Tarsipes spenserae*, of southwestern Western Australia.

Superfamily Vombatoidea

Terrestrial and arboreal Australian marsupials, syndactylism of the second and third digits of the hindfeet. Two extant families; 1 extinct family.

†**Family Diprotodontidae.** Large herbivorous marsupial represented in the Australian fossil record as far back as the Miocene.

Family Phascolarctidae (koalas). Small bearlike Australian marsupials. Adapted to arboreal living and restricted to a diet of the leaves of a few species of eucalyptus. First and second digits of the forelimbs are opposable to the other digits. There is no tail. Pouch opens backward. One species along the eastern Australian coast: *Phascolarctos cinereus*.

Family Vombatidae (wombats). Woodchuck-like marsupials distinct from all others in having a single pair of ever-growing upper and lower incisors and in having rootless, high-crowned cheek teeth. The tail is very short. Includes *Phascolonus*, the giant wombat of the Pleistocene; *Vombatus* (common, or naked-nosed, wombat); *Lastorhinus* (hairy-nosed wombat). Comprises 2 genera and 4 species.

Critical appraisal. An earlier classification, which is still encountered (see section *The class Mammalia*), attributes ordinal rank to the marsupials. Under that system are listed 6 superfamilies and from 13 to 18 families, according to the authority followed. Various attempts to subdivide the order into suborders to accommodate necessary changes required by accumulated taxonomic findings have not been widely accepted. Along with a resurgence of interest in marsupial relationships has developed a growing body of opinion that the order rests too heavily on the presence of a pouch or its vestige. Such a foundation is as unwarranted as the reduction of all recognized orders of placental mammals into a single order because they all possess a placenta. In fact, as many variations and as great a degree of variations are found within the marsupials as are found in the placental mammals. Current opinion, therefore, favours the recognition of the Marsupialia as a superorder comprising four orders as given above.

(H.M.V.D.)

Insectivora (shrews, moles, hedgehogs, tenrecs, solenodons)

The mammal order Insectivora includes the shrews, moles, hedgehogs, and several lesser known groups. Of about 400 species in the order, nearly 300 are shrews of the family Soricidae. All insectivores are small, as mammals go, the largest being about the size of a small rabbit.

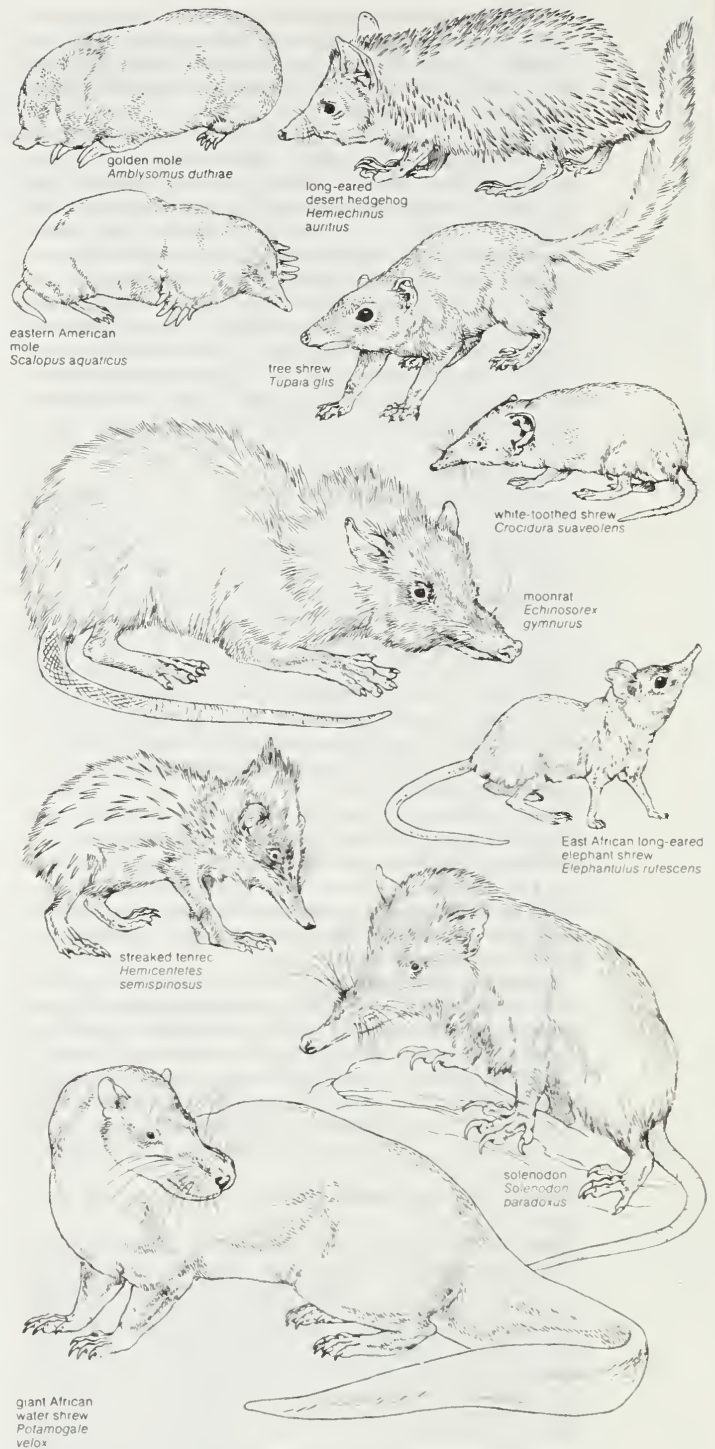


Figure 14: Representative insectivores.

Drawing by R. Keane based on photos courtesy of (golden mole) Herbert Lang through J. Meester, (eastern American mole) U.S. Fish and Wildlife Service (white-toothed shrew) Liselotte Dorfmueller (streaked tenrec) Howard E. Uible (solenodon) New York Zoological Society from (long-eared desert hedgehog) tree shrew, East African long-eared elephant shrew) E.P. Walker, *Mammals of the World*, (giant African water shrew) Grzimek's Tierleben *Enzyklopadie des Tierreiches*, Kindler Verlag Zurich

Shrews are probably the smallest of all mammals; some species of *Sorex* and *Crocidura* weigh as little as 2.5 grams (0.09 ounce). The general body plan (Figure 14) varies greatly within the order. Many shrews are superficially mouselike but may be distinguished from mice and other rodents by the absence of chisel-like incisor teeth. Typical moles (Talpidae) and golden moles (Chrysochloridae) are cylindrical, burrowing animals with short, thick limbs and reduced tails and eyes. Hedgehogs (Erinaceidae) and some hedgehog-like tenrecs (Tenrecidae) are stocky animals whose upper surfaces are covered with short spines. Elephant shrews (Macroscelididae) are rather delicate animals, with long hindlegs like those of kangaroo rats or

jerboas. Within the larger families the body plan is quite variable. Most insectivores live at or near ground level, except for the squirrel-like tree shrews and several aquatic moles, shrews, and tenrecs.

Importance to man. Small and secretive, most insectivores are rarely seen by man and are of almost no economic importance. The velvety fur has poor durability, but mole skins have occasionally been used for trimming various garments and for complete coats.

The burrows of moles in lawns may cause the collapse of the turf, and those in fields and gardens provide small rodents with access to plant roots. Indirectly, the more abundant insectivores are ecologically important as consumers of invertebrates and occasionally of small vertebrates and as prey for larger animals. None is an important disease carrier, although elephant shrews harbor malaria, and some sorcid shrews may carry plague.

A few insectivores are the subjects of aboriginal myths. Some Eskimos believe that the red-toothed shrews (*Sorex*) will attack a man and kill him by burrowing through his body to the heart. The hero shrew (*Scutisorex congicus*), whose backbone of fused and arched vertebrae protects it from crushing, is venerated by central African natives, who believe that consumption of its heart will confer the animal's strength upon the eater.

Distribution. Insectivores are widely distributed but are not found in Antarctica, the Australian region, and South America (except for extreme northern regions). Various marsupials found in Australia are ecologically equivalent to certain insectivores, including a marsupial mole. Why insectivores are absent from South America has not yet been satisfactorily explained.

NATURAL HISTORY

Life cycle. **Reproduction.** Temperate zone insectivores, such as shrews, moles, and hedgehogs, have cyclic reproductive periods. Males develop enlarged testes in late winter or early spring that begin to produce sperm. At the same season females develop enlarged ovarian follicles in preparation for ovulation (release of eggs). Because female hedgehogs may undergo several periods of heat during which copulation does not result in pregnancy, it seems likely that ovulation occurs at a specific time, regardless of other reproductive activities. In shrews and moles, by contrast, every mating results in pregnancy; it has thus been supposed that ovulation in these animals is induced by copulation. An inverse relationship seems to exist between gametogenesis (the production of eggs or sperm) and the development of certain cells in the ovaries and testes. Female moles also exhibit a vaginal cycle, in which the vagina opens to the exterior only during the season of heat, pregnancy, and lactation. The majority of tropical insectivores appear to breed throughout the year; presumably, sperm production is continuous, and ovulation can occur at any appropriate time. Some evidence indicates that, in tropical areas with wet and dry seasons, the majority of births take place with the approach of the rainy season and the accompanying increase in the abundance of invertebrates. In West Malaysia, for example, 41 percent of the recorded pregnancies of one species of white-toothed, or musk, shrew (*Suncus*) occur from October through December, and the peak of the monsoons occurs in January; 19 percent of the pregnancies occur from April through June, before the height of the dry season in July. South African golden moles breed during rainy seasons. Some tenrecs become dormant during the Madagascan dry season, so it is likely that their reproductive period is also related to approaching rains.

Courtship and mating have seldom been observed among insectivores. Males and females of most species remain together for a brief period during which mating takes place. Among tenrecs the male may bite the female's neck; copulation then takes place, with the male mounting the female from behind in all known cases. Shrews may remain locked together for up to five minutes after the male dismounts, while the female drags the male after her. The ejaculate of the male short-tailed shrew (*Blarina brevicauda*) includes a waxy plug that blocks the vaginal orifice, preventing the loss of semen.

The length of the gestation period varies from two or three weeks in shrews to 50 to 60 days in elephant shrews and tree shrews. In part, the length depends upon the condition of the young at birth; shrews and moles, for example, with a gestation of three or four weeks, have naked, blind, and helpless young. Hedgehogs, with a slightly longer period, have young that, although essentially helpless, do have partially developed spines. Tenrecs may have a 50-day gestation; the young have spines and rather quickly develop neuromuscular integration sufficient to enable them at about four days of age to walk and to exhibit at five days the stereotyped tenrec "bucking" defensive behaviour pattern (see below). Newly born elephant shrews have fur, open eyes and ears, and the ability to follow the parent. That a long gestation does not necessarily produce more developed young in insectivores, however, is shown by tree shrews, which bear helpless young after a relatively long gestation period. After giving birth, females of some species, such as tree shrews and some red-toothed shrews, may immediately mate again.

The number of young produced depends in part on species, latitude, and condition of young at birth. Elephant shrews, tree shrews, otter shrews (Potamogalidae), solenodons (Solenodontidae), and golden moles have one to three young. In the hedgehog family the spiny species have one to seven young; the nonspiny tropical Asian gymnures commonly have two. Moles average two to five young, depending on the species. Shrews have more young in temperate areas than in the tropics. Average numbers recorded range from two to three in Africa; two on the island of Guam; three, four, and five in various parts of India; and seven to eight in Great Britain. Various tenrecs may hold the record for number of young among placental mammals, *Centetes* having as many as 25; the presence of more than 30 embryos has been recorded.

Maturation and longevity. Embryonic development has been studied in only a few insectivores. In most shrews, moles, hedgehogs, and tenrecs, the eyes and ears open in about a week, and the babies are weaned in three or four weeks. Tropical white-toothed shrews reach sexual maturity in about five weeks, as do elephant shrews; tenrecs attain this condition in two months, a time that may also apply to other tropical insectivores. Temperate zone shrews, moles, and hedgehogs attain reproductive maturity in the second year of life. The characteristic chain-forming behaviour of young white-toothed shrews (in which the young form a chain behind the mother, each animal gripping with its teeth the rump fur of the one ahead) may begin after the ears open but before the eyelids part, suggesting that this behaviour is organized by sound stimuli.

The average life-span of most wild insectivores, like that of most small mammals, is probably quite short. It has been estimated at six weeks for shrews. Some individuals of populations of temperate zone shrews and moles, however, must live at least one year in order to participate in reproduction. These species seemingly are programmed for a year of healthy life, but few live beyond this; they succumb from wornout teeth. Captive shrews have survived two and one half years, elephant shrews three years, and hedgehogs 10 years.

Response to adverse conditions. As a normal part of their annual cycle, various insectivores have developed strategies for surviving intervals during which food is in short supply. In temperate regions this occurs in winter; in tropical areas dry seasons often are times of privation. Hedgehogs cope with the problem of food shortage by seeking a shelter, rolling themselves into a ball, and becoming dormant. Although the temperature of the extremities drops to that of the surrounding air, that of the heart remains at the level of an active animal; in this way energy is used slowly. In the dormant hedgehog, white blood cells concentrate in unusual numbers around the stomach and intestine, the sugar level in the blood drops by one half, and the level of magnesium in the blood rises to twice the normal value. In a normally dormant hedgehog, the Islets of Langerhans, which are tissues in the pancreas that secrete the hormone insulin, are active. Certain tenrecs are the only other insectivores definitely

Gestation periods

Role of shrews in mythology

Dormancy in hedgehogs

known to become dormant, a state they may enter during the dry season or in cool weather. The body temperature drops to that of the environment, and the respiratory rate may slacken to five or six breaths per minute. One kind of golden mole, *Chrysoptax*, has been reported to become dormant.

Another common stratagem used by many animals to survive lean times, the storage of food, is known for only a few insectivores. Some kinds of shrews store such invertebrates as snails. The European mole (*Talpa europaea*) may paralyze earthworms before burying them in groups in its mound. Energy may be stored as fat by some hedgehogs, tenrecs (*Microgale*), and star-nosed moles (*Condylura*). The last two genera use the tail as a storage organ.

A curious and unexplained aspect of the annual cycle of some temperate zone shrews involves the volume of the brain and the height of the braincase; both decrease to a minimum in midwinter, then increase again.

Behaviour. Feeding behaviour. All insectivores feed wholly or partly on animal matter and spend a substantial amount of time foraging. Shrews consume large quantities of invertebrates and occasionally small vertebrates. Moles eat soil invertebrates, including worms. Larger species eat larger prey; hedgehogs, for example, may consume frogs, mice, birds, lizards, and snakes and in captivity will eat bread and milk. The diet of golden moles is not well-known, although some species are known to prefer legless lizards (*Anguis*); others consume earthworms and insect pupae and larvae. Tree shrews are omnivorous, eating many insects as well as fruit; in captivity they thrive on commercial monkey food. Elephant shrews avidly eat locusts and, in captivity, thrive on bird meat. The aquatic tenrec *Limnogale* eats the tubers of the pondweed *Aponogeton*, and also takes fish. The water shrews consume aquatic animals, sometimes even small fish, and the otter shrew eats crabs, fish, and amphibians.

Many insectivores probe through leaf litter or soft soil. Solenodons proceed in this way until the long snout contacts prey; then the forepaws sweep forward and inward, gathering the object into the mouth. So stereotyped is this behaviour that even in captivity the animal captures inert food items in this way. Young solenodons follow the foraging parent and may learn the scent of acceptable food by smelling the mouth of the female. Young tenrecs may also learn to discriminate food by following their mother. Tenrecs, which often forage in groups, become excited at the scent of earthworms and begin to root enthusiastically. The sound of rooting and chewing stimulates other tenrecs in their searching activity. Shrews search through leaf litter or underground passageways of other mammals, locating prey by smell and contact. Moles do much of their hunting in their burrows, which may act as natural traps into which invertebrates fall. Golden moles emerge after rains and forage in the surface soil. Some shrews hold food items in their paws while eating and may carry food to preferred sites to eat it. Tenrecs may use the forefeet to hold down worms while they are torn in pieces.

A number of insectivores are reported to drink water, although elephant shrews probably do not do so. Tree shrews use water for drinking and bathing. Many insectivores may obtain sufficient water from their food or from dew.

Orientation and activity periods. Except for tree shrews and elephant shrews, both of which are primarily diurnal (*i.e.*, active by day), insectivores rely relatively little on vision for orientation. In most species, olfactory, tactile, and auditory cues are most important. Shrews, some tenrecs, and possibly the solenodons use a crude form of echolocation, a process by which objects are located by sound waves reflected back to the animal. Tenrecs produce clicks with the tongue at frequencies from five to seven kilohertz. Spiny tenrecs produce high frequency noises by vibrating together rows of overlapping spines on the middle of the back; the sounds produced in this way probably are used in social communication rather than orientation. Shrews and solenodons produce, probably with the larynx, clicks of 10 to 31 kilohertz (in solenodons) and 25 to 60 kilohertz (in shrews). Many insectivores, when exploring strange surroundings, maintain contact with a solid sur-

face such as a wall or a log, which is probably tested with the facial sensory hairs (vibrissae). After a territory has become familiar, many insectivores are able to traverse it by memory and may become severely disoriented if the landscape changes.

Most insectivores are active at night, although shrews may be active throughout the 24-hour period. Most elephant shrews and tree shrews are diurnal, but at least one genus in each family is nocturnal. True moles and apparently also the little-studied golden moles are active by day or night.

Locomotion. In walking, most insectivores proceed by moving the right forefoot and left hindfoot, then left forefoot and right hindfoot forward together (crossed extension limb synchrony). Solenodons run with a quadrupedal ricochet (springing) gait, in which both forefeet and then both hindfeet strike the ground together. Other large insectivores also probably run in this way. Elephant shrews may hop bipedally, but their locomotion has not yet been studied. Tree shrews are skillful climbers but spend much of their time on the ground; some tenrecs are also partly arboreal. A few insectivores, such as otter shrews, several species of shrews and moles, and one tenrec, are partially aquatic and swim with alternate strokes of the hindfeet. The American water shrew (*Sorex palustris*) has been observed to scamper over the surface of still water, apparently being supported by the surface tension and the fringe of stiff hairs on its feet. The water tenrec (*Limnogale*) has webbed feet, as do the aquatic moles *Desmana* and *Galemys*. Some aquatic insectivores may use a laterally compressed tail in swimming, such as is found in the water shrew *Nectogale elegans*, the water moles *Desmana moschata* and *Galemys pyrenaicus*, and the otter shrew *Potamogale velox*.

Grooming. All insectivores spend some time in self-grooming. The universality of this behaviour suggests that it is necessary for keeping the hair in proper condition for its role in thermoregulation. Solenodons use only the hindfeet in grooming, perhaps because the forefeet are enlarged for digging, and the elongated snout hinders the use of the mouth in grooming. Shrews use the hindfeet and also the tongue. (Why the forefeet are not used by shrews is not clear.) Hedgehogs likewise do not use the front feet for cleaning. The tenrec *Echinops* uses both hindfeet and front feet. Tree shrews of the genus *Tupaia* use all four feet, but the forepaws and mouth are the dominant cleaning organs. These animals also comb the fur with the lower incisors.

Burrowing and nesting. Probably the most complex insectivore shelters are built by moles. A mole may have an underground nest chamber surrounded by concentric rings of tunnels interconnected by radiating ones. Shallow surface tunnels extend from this deeper complex and are evidenced by raised ridges, which mark their course. The mole may also throw up a large mound of earth on the surface that marks the location of its deep tunnel system. Burrows are dug by using the spadelike forefeet in a manner resembling the breaststroke of swimmers. The golden moles may use the same method, but little is known concerning their excavations. Some shrews dig tunnels, but many species use the runs of other animals. Shrews build surface nests, lined with shredded plant material, in which they rest or bear their young. Solenodons construct nest chambers during the breeding season, as do hedgehogs and tenrecs. The last two also build nests for keeping warm when the temperature drops. Elephant shrews often seek shelter among rocks and probably seldom excavate their own burrows; they may, however, occupy those abandoned by rodents. Tree shrews build nests of leaves and other vegetation, often among roots or fallen timber or in a tree cavity above the ground. When frightened, many insectivores immediately flee to a nearby burrow entrance.

Defensive behaviour. Spiny insectivores, such as hedgehogs and tenrecs, defend themselves by rolling into a ball with the spineless undersurface curled inward. Hedgehogs have a specialized muscle that is capable of pulling the spiny dorsal skin about the animal so that only a small unprotected area remains. Spiny tenrecs, when disturbed, exhibit a stereotyped defense behaviour in which the

Movements in water

Sound production

spines on the head and neck are erected, and the animal "bucks" by throwing its forequarters upward. The action would drive spines into an attacker. Solenodons and tenrecs gape and hiss at intruders, a pattern also seen in the unrelated American opossum, a marsupial. Shrews, when confronted with strangers, open the mouth and emit bursts of high-pitched cries.

Shrews and solenodons have been shown to produce saliva with poisonous properties; it probably functions in quieting large prey animals. Shrew bites may be painful to humans.

Social behaviour. The majority of insectivores are rather antisocial, except at breeding time and in the context of the female-young relationship. Adult shrews and moles tolerate other adults of their species briefly during the period of mating, following which the female and young form a group until the young are full grown; the young rarely appear alone during this interval. Young white-toothed shrews (*Crocidura*) exhibit chain formation, in which each animal grasps the rump fur of another, one grasping the mother, after whom they then trail. Baby shrews are rather indiscriminate and may attempt chain formation with such inappropriate subjects as white mice.

Juvenile
behaviour

Solenodons have a prolonged period of juvenile dependency during which the young follow the female, perhaps learning to locate favoured foraging areas and to discriminate in food selection. Adult solenodons approach each other with open mouths, perhaps emitting high-frequency clicks. Contact involves one animal closing its mouth over the snout of the other, then pressing the snout against its own flanks and back. Several solenodons may occupy a burrow. Various tenrecs forage in large groups which probably consist of one or more females and young. Elephant and tree shrews live solitarily or in groups of two or three, consisting of a female and her young. Hedgehogs may also form temporary groups of female and young.

Hedgehogs and spiny tenrecs perform an unexplained "self-anointing" behaviour, which is triggered by the presence of such substances as urine or certain chemicals. The animal licks the material until a frothy spittle is produced; in some species this is then spread over the body by the tongue in typical hedgehogs (*Erinaceus*) or by the mouth and forefeet in tenrecs.

Ecology. *Predation upon insectivores.* Shrews are preyed upon by a variety of birds, mammals, and snakes, but seemingly the strong odour produced by the flank gland of shrews makes them unpalatable to some predators. Skulls of shrews are commonly found in owl pellets (regurgitated indigestible remains of prey). Little is known about the predators of elephant and tree shrews. True moles and golden moles may be safe from many predators because of their subterranean habits; spiny hedgehogs and tenrecs may, for the same reason, be relatively safe. In Madagascar some tenrecs are preyed upon by man. Solenodons seemingly suffered little predation on Cuba and Hispaniola until man introduced dogs, cats, and mongooses there. These carnivores have drastically reduced solenodon populations, probably because the solenodons had evolved no defenses against such predation.

Territoriality. Each adult shrew occupies its own territory of about 0.1 to 0.4 hectare (0.25 to one acre). About 40 percent of its behaviour involves avoiding other shrews. When shrews from adjacent territories encounter one another, they engage in mutual avoidance if possible. Avoidance is made easier by an oily odoriferous substance exuded by the flank glands, which helps to warn off other shrews of the same species. A similar kind of territoriality probably is practiced by moles. Population structure of most other insectivores has not yet been studied. Tree shrews apparently defend nest territories, and males may mark these areas by smearing onto branches the yellow exudate from throat and chest glands.

Habitats. Insectivores utilize diverse habitats. Shrews occupy the ground surface in forests and grasslands but are less common in deserts, although *Notiosorex* and *Diplomesodon* occupy New and Old World deserts, respectively. Moles are largely confined to North Temperate Zone forests and meadows, especially in arcas of deciduous forests and adjacent prairies. Hedgehogs live in forests,

grasslands, and deserts of the Old World. Spineless hedgehogs (Echinosoricinae) live in tropical rain forests of Asia. Golden moles occupy many habitats in southern Africa. Otter shrews live in riparian situations (*i.e.*, on the banks of natural water courses), in the Congo Basin and environs. Elephant shrews prefer African savanna areas, but some are forest dwellers. The tenrecs live in most habitats on Madagascar; some are partly arboreal. Tree shrews are largely terrestrial in tropical Asian forests, but most can climb well and are partly arboreal. Solenodons are strictly terrestrial foragers in the forests of Cuba and Hispaniola.

FORM AND FUNCTION

The order Insectivora is the most difficult of living mammalian orders to define because the groups included retain traits of primitive placental mammals and are not very specialized. Although each family has specializations, they are not such as to suggest that each group should be raised to ordinal level. Furthermore, many of the families are as old as many mammalian orders, and the evidences of common relationship are often obscured.

Skull. The skull of insectivores is typically low in that the braincase does not rise abruptly from the level of the rostrum (the bony support for the snout), which tends to be long and tapered. Seen from above, the skull is triangular in outline. The cerebral cortex of the brain is not well developed, and the olfactory bulbs are large, in conjunction with the well-developed sense of smell. In addition, the elongated rostrum supports muscles that activate a still longer, highly mobile snout with well-developed tactile capabilities that are useful in probing for food. The middle ear cavity is partly or wholly enclosed by the tympanic bone. In such forms as shrews and hedgehogs, the tympanic bone (surrounding the middle ear) is ringlike and attached to the base of the skull (basicranium) only at one point, probably the primitive situation. Other forms, such as moles and golden moles, have an expanded tympanic bone fused to the basicranium to form an auditory bulla (an enclosed chamber). The bulla of tree shrews and elephant shrews includes an entotympanic bone, a trait also seen in primates. Many insectivores have small external ear pinnae (the projecting "shells" of the ear); this is nearly or quite lacking in burrowing forms such as moles but may be rather prominent in such forms as desert hedgehogs. In shrews and tenrecs, at least, sensitivity to high-frequency sounds is well developed.

Presence
of
tympanic
bone

Orbital fossae (eye sockets) and eyes are variably developed in insectivores, most species having relatively small eyes. In moles and golden moles, in which the eyelids are fused shut, probably only differences in light intensity can be detected. Large eyes and well-developed vision are present only in tree and elephant shrews. In the former the orbital fossa is completely encircled with bone.

Dentition. Insectivores have from 26 to 44 teeth, the latter being the maximum for placental mammals. Incisors range from three to six pairs. The second lower incisors of solenodons are long and grooved and transmit venom from the submaxillary salivary glands, the ducts of which end at the bases of these teeth. The upper and lower first incisors of shrews are provided with one or more accessory cusps. These specialized teeth close so as to provide the animal with a delicate tool for picking up small invertebrates. The incisors of tenrecs may also be provided with accessory cusps, but tenrecs lack the other dental specializations of shrews. The lower incisors of tree shrews form a comblike structure and are used for grooming the fur. The canine teeth, though present in insectivores, are usually small, rarely extending much above the level of the other teeth; they are not strong and elongated as those in carnivores (an exception is the Palearctic mole *Talpa*). The canine teeth of hedgehogs and moles have double roots, an unusual condition among mammals. Premolars range from the simple one-cusped structures seen in shrews and some moles to multicusped molariform teeth. The molariform premolar is usually the fourth or, more rarely, the third.

Insectivoran molars are of a primitive placental type. The crown supports several usually conical cusps arranged in a triangular pattern (called protocone, paracone, and

Molars

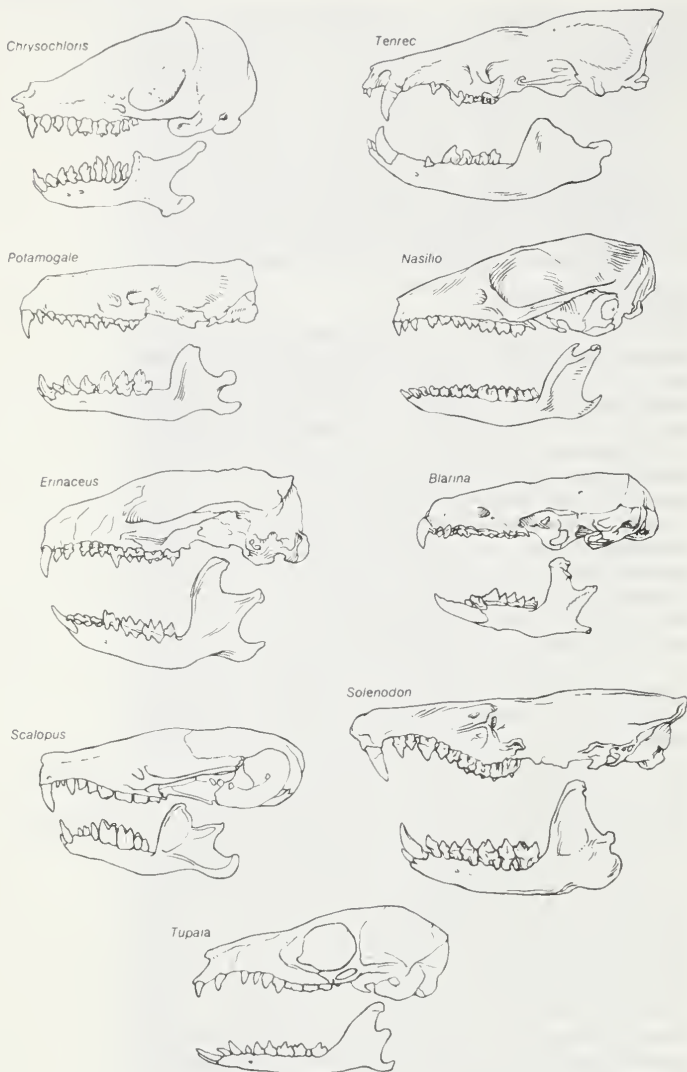


Figure 15: Skulls (cranium and lower mandible) of representative insectivores.

Drawing by R. Keane based on (Chrysochloris, Tenrec, Potamogale, Nasillo) P.-P. Grasse *Traite de Zoologie*, vol. 17 (©1955), Masson & Cie, Paris. (Blarina, Scalopus) E. R. Hall and K. P. Nelson. *The Mammals of North America*, copyright © 1959. The Ronald Press Company, New York. (Tupia) *Proceedings of the U.S. National Museum*, vol. 45.

metacone on the upper teeth; protoconid, paraconid, and metaconid on the lower). The upper triangle is often squared off by the addition of a posteromedial hypocone (*i.e.*, on the rear inner corner). To the lower triangles is added a posterior shelf (talonid) frequently made up of cusps called hypoconid and entoconid. These molars close so as to allow some degree of both slicing and crushing action; the term tuberculosectorial is applied to them. Omnivorous forms tend to emphasize the squaring up of the teeth and the crushing, rather than the slicing, function by lowering the cusps so that the closing surface is made up of low hillocks. This type, the bunodont molar, is best developed in hedgehogs. When the triangular part of each crown is emphasized and elevated, the result generally indicates a carnivorous diet; such teeth are seen in tenrecs, otter shrews, golden moles, and solenodons. In these four groups the paracone and metacone (outer cusps, forming the base of the upper triangle) are close together and medially removed from the outer edge of the tooth. Such a molar type is called zalambdodont. Moles and shrews have converted the paracone and metacone into V-shaped ridges, such that the two produce a W-shaped pattern on the upper molar (termed dilambdodont). The result is that there are four, rather than two, cutting edges on the upper molar. This specialization, seen also in insectivorous bats, is seemingly especially useful in chopping up small invertebrates. The molars of tree shrews and elephant shrews are also dilambdodont; in the former the appearance is sufficiently similar to that of true shrews to suggest some

dependence on insects. The teeth of elephant shrews show a curious parallel to those of artiodactyls in being somewhat high crowned, with considerable molarification of the premolars. The extent to which elephant shrews are vegetarians, as their dentition suggests they may be, has not been measured. Shrews may replace their milk teeth while still in the uterus.

Limbs. The appendicular skeleton of insectivores retains many traits seen in primitive placental mammals. All except the otter shrews retain a clavicle connecting scapula (dorsal part of pectoral girdle) to sternum (breastbone). Since absence or reduction of the clavicle usually characterizes mammals that are specialized for rapid leaping locomotion, such as artiodactyls (antelopes, pigs, and cattle) and perissodactyls (horses, rhinoceroses, and tapirs), its absence in otter shrews, which are reported to be clumsy on land, is perplexing. The humerus (upper bone of the forelimb) of insectivores retains an entepicondylar foramen (an opening in the articular surface), absent in many more specialized mammals. The radius and ulna (bones of the forearm) are separate and well developed. The manus (forefoot) is rarely specialized, usually possessing five digits and a relatively full complement of wristbones (carpals). In some forms (shrews and hedgehogs) one wristbone, the centrale, is missing or fused with an adjacent bone, as it is in most primates. Moles have broad forefeet, in which supernumerary bones are closely knit to form a rigid plane, as in the blade of a shovel. Golden moles have the hand digits reduced to four, two of which are greatly enlarged and bear huge, picklike claws. The pelvic girdle shows no peculiarities. In a few families (shrews, moles, golden moles, and otter shrews) the pubic bones do not form a symphysis (union); in others a short or long symphysis is present. The femur (thighbone) is marked by a laterally placed third trochanter (a ridge) for muscular attachment, presumably a primitive trait but one seen in a variety of other mammals. The tibia and fibula (lower bones of the hindlimb) are usually fused distally (near the ankle), but in solenodons this fusion takes place only with old age and in some tenrecs and all tree shrews not at all. Fusion reduces the lateral mobility of the foot but promotes more efficient forward and backward motion. In hedgehogs and elephant shrews the toes may be reduced to four. All insectivores except elephant shrews have plantigrade foot posture; *i.e.*, the entire sole is placed down in walking. In some species, however, the heel is frequently lifted from the ground.

Fusion of tibia and fibula

Tail. The tail structure varies according to the behaviour of the animal. Burrowing insectivores have little or no tail, in common with most other burrowing animals. Climbing and running species have well-developed tails, haired or nearly naked. The slow-moving hedgehogs generally have short tails. In aquatic forms, such as aquatic shrews, moles, and potamogales, lateral compression of the tail seemingly suits it for steering functions.

Fur texture and coloration. Specializations of the hair are seen in hedgehogs and some tenrecs, which possess

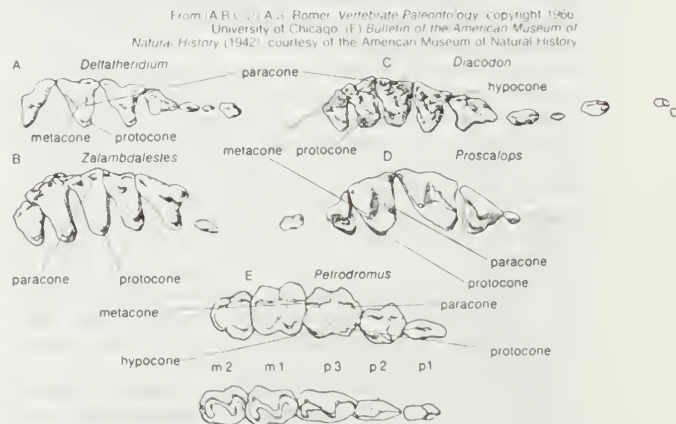


Figure 16: Major molar types in insectivores. (A) An early zalambdodont placental. (B) A Cretaceous leptictid. (C) An early hedgehog with omnivorous bunodont tooth type. (D) A dilambdodont mole. (E) A putatively vegetarian elephant shrew (see text)

spiny hair. Burrowing species have reversible fur, which can lie down either anteriorly or posteriorly as the animal moves forward or backward in the tunnel. The fur of golden moles has a distinctive metallic sheen or iridescence. Few insectivores are strikingly coloured; most are brown, gray, or black with little countershading. Exceptions occur in diurnal species, which may be more brightly coloured and have some indication of patterns.

Digestive and reproductive specializations. Tree shrews and elephant shrews possess a cecum (a blind pouch) at the junction of small and large intestines. An intestinal cecum is not present in carnivorous insectivores, since it serves for storage and continued digestion of plant food. Elephant shrews, whose teeth suggest a diet consisting largely of plants, have a large cecum; tree shrews, which probably eat relatively less plant food, have a comparatively small cecum.

A true scrotum containing the testes is present only in tree shrews; in some other forms, such as moles and tenrecs, the testes lie close to the surface in the perineal region (the area ventral to the anus) but do not leave the abdominal cavity. A penis bone (baculum) is found in tenrecs and moles. Elephant shrews have a trifurcate penis, a trait that led some early students to ally them with the marsupials, some of which have a bifurcate structure. The form of the uterus is Y-shaped, or bicornuate, as in many other mammals.

Metabolism. Although insectivoran physiology has been studied in connection with special abilities, such as dormancy in hedgehogs, it is not well-known for the order. Of great interest are the smaller shrews, which, because of their minute size and relatively large ratio of body surface to volume, have an exceptionally high rate of metabolism, necessitated by the high rate of energy loss in the form of heat. As a result shrew life consists largely of a frenetic search for food. A captive female *Sorex* ate 3.3 times her own weight every 24 hours.

EVOLUTION AND CLASSIFICATION

Evolution. In late Cretaceous and early Cenozoic times (beginning about 100,000,000 years ago) a number of primitive placental mammals existed that did not display the specialized characteristics of modern orders; they possessed tritubercular (three-cusped) cheek teeth, with the paracone and metacone well separated and located laterally on the occlusal (closing) surface. These animals, the Leptictidae, some of which persisted until Oligocene times (about 38,000,000 years ago), were close structurally to Eocene hedgehogs and may well be ancestral not only to many modern insectivores, especially hedgehogs, shrews, and moles, but also to a variety of other modern mammals, such as ungulates and primates. These insectivores, in which the paracone and metacone are well separated and are near the lateral edge of the molar, are thought by some paleontologists to have descended from Late Cretaceous leptictids and to form a separate lineage from the zalambdodonts. The fossil record is fairly clear in showing common ancestry of shrews and moles, and probable derivation from late Eocene or early Oligocene hedgehog-like ancestors (perhaps 40,000,000 years ago). Some paleontologists believe that the presence of zalambdodont molars in tenrecs, otter shrews, and golden moles is indicative of independent descent from the zalambdodont deltatheridiids, a family of doubtful ordinal affinities, of Late Cretaceous time. It has recently been suggested that solenodons also descended from a hedgehog-like ancestor, which reached the Greater Antilles in early Tertiary time. Tree shrews are known from Paleocene time (about 65,000,000 years ago) and probably come closest among living mammals to representing the appearance of the early Cenozoic primitive placentals. Fossil elephant shrews from the African Miocene (about 26,000,000 years ago) offer few clues to the origin of the Macroscelididae, which were in all likelihood derived from the unspecialized Early Cenozoic placental group. Tenrecoids and golden moles have an exclusively African and Malagasy fossil history, also stemming from the Miocene.

Classification. *Distinguishing taxonomic features.* Insectivores are characterized by the lack of distinctive fea-

tures seen in other mammalian orders rather than by the possession of unifying morphological traits. All possess a low braincase, a rather long conical rostrum (snout) and unspecialized legs, feet, and locomotion. They can be distinguished from similarly built rodents by the lack of the gnawing incisors (followed by a diastema, or space) that characterize the rodents and from similarly built marsupials by the lack of an abdominal pouch.

Annotated classification. Taxonomists disagree on the relative degrees of relationship among insectivoran families (except that shrews, moles, and hedgehogs are conceded to form a cohesive phyletic unit); the classification omits subordinal and superfamilial groupings.

ORDER INSECTIVORA

Primitive placental mammals characterized by a low braincase, conical snout, and lacking locomotor or dental specializations of other mammalian orders; 77 genera, 406 species.

Family Erinaceidae (hedgehogs and gymnures)

Eocene to Recent; Eurasia and Africa. Spiny pelage (covering) present (subfamily Erinaceinae, hedgehogs) or lacking (Echinosoricinae, gymnures); zygomatic arch (below eye orbit) present in skull; molars quadrate, bunodont (see above *Dentition*); auditory bulla a tympanic ring, not fused to cranium. Terrestrial, omnivorous, often hibernating or estivating. Hedgehogs are Eurasian and African; gymnures native to the Asian tropics, often shrewlike. Ten genera, 14 species.

Family Talpidae (moles)

Eocene to Recent; Holarctic. Manus very broad and specialized for digging; molars with W-shaped crests; zygomatic arch present; auditory bulla fused to cranium; eyes minute; ears without pinnae; humerus bone extremely broad and short. Burrowing (fossorial), sometimes aquatic. Fur of some commercial value. Fifteen genera, about 22 species.

Family Soricidae (shrews)

Eocene to Recent; Africa, Eurasia, North America, northern edge of South America. First upper and lower incisors procumbent, multicusped; zygomatic arch lacking; bulla a tympanic ring, unfused to skull; molars with W-shaped crests; small size. Terrestrial, aquatic, or subfossorial foragers in ground litter or water for small invertebrates. Twenty-four genera, approximately 291 species.

Family Solenodontidae (solenodon, almiqui)

Pleistocene to Recent; Cuba and Haiti. Molars zalambdodont; zygomatic arch incomplete; incisors canine-like, lower ones grooved; long scaled tail; long mobile snout. Terrestrial, perhaps nearing extinction. Two genera, 2 species.

Family Tenrecidae (tenrecs)

Miocene to Recent; Madagascar. Molars zalambdodont; zygomatic arch incomplete; bulla composed of a tympanic ring and basicranial elements; form variable from nearly tailless, spiny, rabbit-sized *Tenrec* to long-tailed, soft-furred, shrewlike *Microgale*. *Oryzorictes* is molelike, while *Linnogale* is muskrat-like and a vegetarian. Nine genera, 20 species.

Family Potamogalidae (otter shrews)

Recent; West Africa. Webbed feet, otter-like form; zalambdodont molars. Riparian, aquatic, feeding on aquatic vertebrates and invertebrates. Two genera, 3 species, often included as a subfamily of Tenrecidae.

Family Chrysochloridae (golden moles)

Miocene to Recent; Africa. Two digits on front feet with enormously enlarged claws; proximal segments of arm recessed into sides of thorax; heart and lungs elongated; zalambdodont molars; complete zygomatic arch composed of maxillary rather than jugal bone. Eyes covered with skin; tail rudimentary; pelage with metallic sheen or iridescence. Fossorial. Five genera, 11 species.

Family Macroscelididae (elephant shrews)

Miocene to Recent, Africa. Cecum well-developed; feet greatly elongated; molars quadrate and high-crowned; eyes large; auditory bullae complete, including entotympanic bone; zygomatic arch complete; resembling long-snouted kangaroo rats or jerboas. Diurnal, terrestrial, saltatory, animalivorous and phytophagous. Five genera, about 28 species.

Family Tupaiidae (tree shrews)

Paleocene to Recent, tropical southeast Asian mainland, Sumatra, Borneo, Philippines. Small cecum; orbital fossa ringed with bone; zygomatic arch complete; molars quadrate with W-shaped crests; eyes large; auditory bulla includes entotympanic bone; generally like a long-nosed squirrel in appearance and behaviour. Terrestrial and arboreal; diurnal or nocturnal. Five genera, 15 species.

In addition to the above, 14 families of primitive extinct mammals are sometimes assigned to the Insectivora. These families, with their geologic time spans, are listed below.

- Family Dimylidae. Oligocene–Palcocene.
- Family Picrodontidae. Paleocene.
- Family Zalambdalestidae. Cretaceous.
- Family Leptictidae. Cretaceous–Oligocene.
- Family Adapisoricidae. Paleocenc–Miocene.
- Family Pantolestidae. Paleocene–Oligocene.
- Family Apternodontidae. Eocene–Oligocene.
- Family Apatemyidae. Paleocene–Oligocene.
- Family Mixodectidae. Paleocene.
- Family Anagalidae. Eocene–Oligocene.
- Family Paroxyclaenidae. Eocene–Oligocene.
- Family Ptolemaiidae. Oligocene.
- Family Pentacodontidae. Paleocene–Eocene.
- Family Plesiosoricidae. Eocene–Pliocene.

Critical appraisal. It is possible, though current opinion does not support it, that insectivores with zalambdodont teeth (tenrecs, solenodons, golden moles) form a monophyletic unit descended from either pre-tritubercular Mesozoic (Jurassic, about 190,000,000 years ago) pantothers or the zalambdodont deltatheridiids of late Cretaceous and Paleocene times. More likely, zalambdodont teeth were derived from the tritubercular type seen in leptictids. Zalambdodont and dilambdodont teeth represent different ways of increasing the length of shearing edge of a tritubercular tooth. Tupaiids and macroselidids share a number of traits that have led some authorities to separate them as a group called Menotyphla, while grouping the other insectivores as the Lipotyphla. Some authorities have considered the colugos (Cynocephalidae) to belong in the Insectivora, but most now place them in a separate order, Dermoptera. Tupaiids share some traits with primitive primates and have been included in that order by some prominent authorities (see, for instance, the section *Primates*). Elephant shrews have many of these same traits, however, and would have to be included in the Primates as well. Some mammalogists treat the elephant shrews as a distinct order, Macroselidea (see section *Mammalia*). A growing body of evidence suggests that soricids, talpids, erinaceids, solenodons, tenrecs, otter shrews, and possibly chrysochlorids form a naturally related group to which the ordinal or subordinal term Lipotyphla should apply. This group, as well as elephant shrews, tree shrews, primates, bats, and others, apparently have been separate from one another since earliest Tertiary times, and consistency may in the end dictate assignment of a separate ordinal name to each. (J.Fi.)

The Lipotyphla

Chiroptera (bats)

Bats, which comprise the mammalian order Chiroptera, are the only mammals to have evolved true flight. This ability, coupled with the benefits deriving from their system of acoustic orientation (so-called bat sonar), has made the group a successful one in numbers of species and individuals. About 900 species are currently recognized, belonging to some 174 genera. Many species are enormously abundant. Observers have concluded, for example, that some 100,000,000 female Mexican free-tailed bats (*Tadarida brasiliensis mexicana*) form summer nursery colonies in Texas, where they produce about 100,000,000 young in five large caves. Adult males of this species, although equal to the females in numbers, may not range as far north as Texas. (Individuals of the species also range widely throughout tropical America.) Thus, bats of one species alone number at the very least in the hundreds of millions of individuals.

GENERAL FEATURES

Most bats are insectivorous. Little is known of the spectrum of insect species consumed, but the quantities are formidable. The Mexican free-tailed bats of Texas have been estimated to consume about 20,000 tons of insects per year. Bats would thus seem to be important in the balance of insect populations and possibly in the control of insect pests. Some bats feed on fruit and aid in dispersing seeds; others feed on pollen and nectar and are the principal or exclusive pollinators of a number of trop-

ical and subtropical plants. The true vampires of tropical America feed on the blood of large birds and mammals, occasionally becoming significant pests of livestock and sometimes serving as carriers of rabies.

Certain aspects of the physiology of some bats, particularly those involving adaptations for long hibernation, daily lethargy, complex temperature regulation, acoustical orientation, and long-distance migrations, are of interest to experimental scientists.

In tropical countries, in particular, large colonies of bats often inhabit houses and public buildings, attracting attention by their noisiness, guano (droppings), and collective odour. In the West, bats have been the subject of unfavourable myths; in parts of the Orient, however, these animals serve as symbols of good luck, long life, and happiness.

Noxious bats

Diversity of structure. All bats have a generally similar appearance in flight, dominated by the expanse of the wings, but they vary considerably in size. The order is usually divided into two well-defined suborders: the Megachiroptera (Old World) and the Microchiroptera (worldwide). Among members of the Megachiroptera, a flying fox, *Pteropus vampyrus*, may have a wingspread of about 1.5 metres (about five feet) and a weight of about one kilogram (2.2 pounds). The largest insectivorous bat is probably *Cheiromeles torquata*; it weighs about 250 grams (about nine ounces). The largest of the carnivorous bats (and the largest bat in the New World) is *Vampyrum spectrum*, with a wingspread of over 60 centimetres (24 inches). The tiny Philippine bamboo bat, *Tylonycteris pachypus meyeri*, has a wingspread of barely 15 centimetres (six inches) and weighs about 1.5 grams (about 0.05 ounce).

Bats vary in colour and in fur texture. Facial appearance, dominated by the muzzle and ears, varies strikingly with family and often with genus. In several families, a

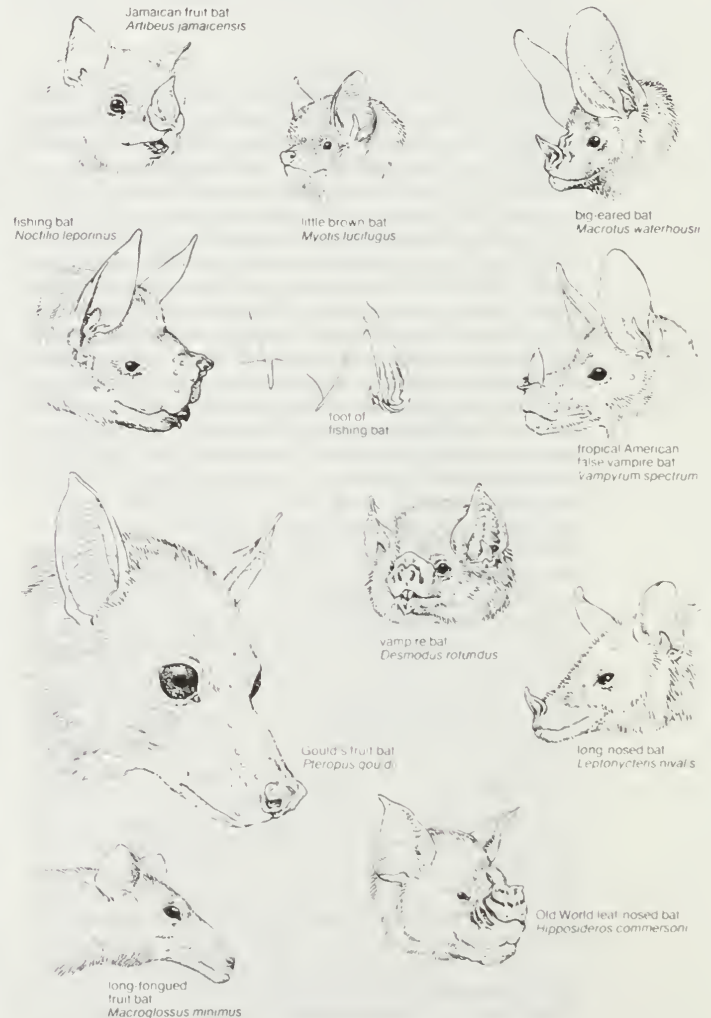


Figure 17: Heads of various bats.

complex fleshy adornment called the nose leaf surrounds the nostrils. Wing proportions are modified according to modes of flight. The tail and the interfemoral (between the legs) membrane also differ, perhaps with feeding, flight, and roosting habits. Finally, at the roost, bats vary in the postures they assume, particularly in whether they hang suspended or rest on the wall and in the manner in which the wings are folded and used.

Distribution. Bats are particularly abundant in the tropics. In West Africa, for example, 31 genera embracing 97 species have been catalogued; in the United States 15 genera, totalling 40 species, are known. Of the 17 living families, three (the Vespertilionidae, Molossidae, and Rhinolophidae) are well represented in the temperate zones. A few phyllostomatid species range into mild temperate regions. Several vespertilionids range well into Canada.

The Vespertilionidae are found worldwide, excepting only the Arctic and sub-Arctic, Antarctic, and isolated islands. The genus *Myotis* has a range almost equal to that of the order. The Molossidae and Emballonuridae also encircle the earth but are restricted to the tropics and subtropics. The Rhinolophidae extend throughout the Old World, the Hipposideridae and Pteropodidae throughout the Old World tropics, and the Phyllostomatidae throughout the New World tropics and slightly beyond. The other families have more restricted ranges.

IMPORTANCE TO MAN

Bats are important to man primarily because they affect other living things through predation, pollination, and seed dispersal. Although losses of commercially valuable fruit are normally small, bananas and figs must, in some cases, be protected by early harvest or by nets. Vampire bats are considered serious pests of livestock in some parts of tropical America, because the small wounds they cause provide egg-laying sites for botflies and because the vampires may transmit rabies and trypanosomiasis to cattle. Other bats also carry rabies or related viruses.

Bat guano as fertilizer The guano deposits of insectivorous bats are used for agricultural fertilizer in many countries and, in the past, were used as a source of nitrogen and phosphorus for munitions. Large guano deposits, in addition, cover and thus preserve many archaeologically interesting artifacts and fossils in caves.

In some parts of Southeast Asia and on some Pacific islands, flying foxes (*Pteropus*) are hunted for food. Small bats are also widely but irregularly eaten.

In species and numbers, bats constitute an important and, on the whole, nonintrusive part of wildlife. Several zoos have established interesting exhibits of bats; indeed, some flying foxes and fruit bats have been exhibited in European zoos since the mid-19th century. Bats are interesting pets and have been kept widely for research purposes but require somewhat specialized care.

NATURAL HISTORY

Life cycle. Details of the life cycle are known for only a few species from North America or Europe. In these, there is an annual cycle of sexual activity, with birth taking place between May and July. In males, the testes, normally located in the abdominal region, descend seasonally into the scrotum, and active spermatogenesis occurs. In females, sexual receptivity may be associated with egg maturation and release. Tropical bats may exhibit a single annual sexual cycle or may be diestrous (*i.e.*, with two periods of fertility) or polyestrous (with many).

The sexual cycles of entire populations are closely synchronized, so that almost all mating activity occurs within a few weeks. The periods of development (gestation), birth, lactation, and fledging are also usually synchronized. Gestation varies in duration: five or six months in *Pteropus*, more than five months in *Desmodus*, three months in some small *Hipposideros*, and six or seven to 14 weeks in several small vespertilionid genera. The length of gestation may be influenced by ambient (surrounding) and body temperature.

In several North American and Palearctic (northern Eurasian) vespertilionids and rhinolophids that hibernate, copulation occurs in the fall, and the sperm are stored in

the female genital tract until spring. Ovulation, fertilization, and implantation occur after emergence from hibernation, when the female again has available an abundant food supply and warm roost. The favourable environmental conditions greatly enhance the young bat's chances of survival.

Most bats bear one young, but the big brown bat (*Eptesicus fuscus*) may bear twins, and the red bat (*Lasiurus borealis*) bears litters of one to four.

At birth the young, who may weigh from one-sixth to one-third as much as the mother, usually have well-developed hindlegs with which they hold on to their mother or to the roost. The wings are very immature. The young, nude or lightly furred, are often briefly blind and deaf. Female bats normally have one pectoral (at the chest) or axillary (at the armpit) mammary gland on each side. Several families that carry their young while foraging also have a pair of false pubic nipples, which the infant may hold in its mouth when its mother flies. The infants are nourished by milk for a period of about five or six weeks in many small Microchiroptera and for five months in *Pteropus giganteus*. By two months of age, most of the Microchiroptera have achieved adult size, having begun to fly and forage three or four weeks earlier.

In many species, females late in pregnancy migrate to special nursery roosts, in which large numbers of pregnant females may aggregate, usually to the exclusion of nonpregnant females, males, and bats of other species. In some cases, the nursery roosts seem to be chosen for their high temperature, which may derive from the sun, the bats themselves, or from decomposing guano. When foraging, some bats (*Erophylla*) leave their infants hanging quietly, one by one, on the cave wall or ceiling. In other cases (*Tadarida brasiliensis*, *Natalus*), the closely spaced infants may move about and mingle on the wall. Among others, especially molossids in buildings, the infants may form play groups and rough-and-tumble on the floor. Some bats carry their young with them. Generally, each mother, on returning to her roost, seeks out her own child by position, smell, and acoustical exchange. In the huge nursery colonies of *Tadarida brasiliensis mexicana* and *Miniopterus schreibersii*, however, the mothers appear to feed the first one or two infants they encounter on returning. This is called the milk-herd system.

Nursery
roosts



T S Lal-PIP

Mother dog-faced fruit bat (*Cynopterus sphinx*), with clinging young, in flight.

Some bats achieve sexual maturity in the first year; others in the second year. Infant mortality appears to be high. Developmental and genetic errors and disease take their toll, but accidents seem to cause more serious losses—the young may fall from the ceiling or have serious collisions in early flight attempts. A fair number of bats probably fail to make the transition from dependent infants to self-sufficient foragers.

Adult bats have low mortality. Predation is rarely serious, especially for cave-dwelling species. Disease, parasitic infestation, starvation, and accidents apparently take small

Life-span

tolls. There are records of several big brown (*Eptesicus fuscus*), little brown (*Myotis lucifugus*), and greater horse-shoe bats (*Rhinolophus ferrum-equinum*) that have lived more than 20 years. Probably many bats in temperate climates live more than 10 years. Longevity has not been established for tropical species.

Several factors probably contribute to the unusual longevity of bats. Generally isolated roosts and nocturnal flight substantially protect them from predation, from some elements of weather, and from exposure to the sun. The largely colonial way of life may ensure that entire populations experience contagious infection and subsequent immunity; indeed, such a pattern may, in the past, have hastened adaptation to disease. The persistent use of various seasonal roosts probably ensures isolation and security, food and water supply, and access to mates. Many bats, moreover, drop their body temperature at rest. Not only is there a probability that this conserves some cellular "machinery," since metabolism is reduced, but fewer hours need be spent in actively seeking food and water.

Behaviour. *Activity patterns.* Nocturnal activity is a major feature of the behavioural pattern of bats: nearly all species roost during the day and forage at night. Carnivorous bats, vampires, and perhaps fishing bats may have an advantage at night over inactive or sleeping prey. In addition, nocturnal flight protects bats from visual predators, exposure to the sun, high ambient temperature, and low relative humidity. The large area of naked wing skin might mean that bats would absorb rather than radiate heat, if they were active during the day. They would also lose body water required for temperature regulation and would then be forced to forage near water or to carry extra water in flight.

The nocturnal-activity pattern in bats is probably kept in synchrony with changing day lengths by their exposure to light at dusk or dawn. Bats often awaken and fly from the cave exit well before nightfall. Should they be too early, their internal clock may be reset. A few species of bats, *Lavia frons* and *Saccopteryx bilineata*, may forage actively during the day, but little is yet known of their special adaptations.

Locomotion. Flight is the primary mode of locomotion in all bats, although the flight styles vary. Some groups (the Molossidae, for example), adapted for flight in open spaces and often at high altitudes, have long, narrow wings, swift flight, and a large radius of turning. Other bats (the Nycteridae, Megadermatidae, and the Glossophaginae), adapted for hovering as they pick prey off vegetation or feed on flowers, have short, broad wings, slow flight, and a small radius of turning. Some bats take flight easily from the ground: members of the genus *Macrotus* do so simply by flapping, vampires (*Desmodus*) leap into the air and then spread their wings and fly. The molossids, however, roost well above the ground since, on takeoff, they fall before becoming airborne.

Though flight speeds in nature are hard to measure, four vespertilionid species, carefully observed, have been timed on the average at 18.7 to 33.3 kilometres (11.7 to 20.8 miles) per hour. In flight, the posture of each of the four fingers incorporated into the wing is under precise and individual control. Finger and arm postures, which determine the shape, extension, and angle of the wings, govern such actions as turning, diving, landing, and hovering. Except when interrupted by insect catches or obstacles, bat flight paths are straight. Insects may be pursued and captured at a rate of up to two per second; during each catch, the flight path is interrupted and thus appears erratic.

In many cases (especially in the families Nycteridae, Megadermatidae, Rhinolophidae, and Hipposideridae), there is little locomotion other than flight. These bats may move across the cave ceiling from which they hang, by shifting their toehold, one foot at a time. A few genera (especially among the Pteropodidae) may crawl along branches, in a slothlike posture, using their thumb claws as well as their feet. The Emballonuridae and Rhinopomatidae hang on vertical surfaces suspended by their hind claws, but with their thumbs and wrists propped against the surface. In this orientation, they can scramble rapidly up or down and forward or backward, as well as sideways.

Bats of many genera (Vespertilionidae, Molossidae, Nocilionidae, Desmodontidae) walk or crawl on either horizontal or vertical surfaces using hindfeet, wrists, and thumbs. Many move freely either backward or forward, a convenience for entering and leaving crevices. The vampires may also leap from roost to roost. The Thyropteridae and Myzopodidae, as well as the vespertilionid *Tytonycteris*, have specialized wrist and sole pads for moving along and roosting on the smooth surface of leaves or bamboo stalks.

Bats are not known to swim in nature except, perhaps, by accident. When they do fall into the water, however, they generally swim competently.

Roosting. Bats choose a variety of diurnal roosts. Each species favours a particular kind of roost, though this varies with sex, season, and reproductive activity. Many bats favour isolated or secure roosts—caves, crevices in cliff faces, the interstices of boulder heaps, tree hollows, animal burrows, culverts, abandoned buildings, portions of buildings inaccessible to man (*i.e.*, roof, attic, hollow wall), and the hollow core of bamboo stalks. Some species roost externally—on tree trunks or in the branches of trees, under palm leaves, in unopened tubular leaves, or on the surface of rocks or buildings. For some, the darkness, stability of temperature and humidity, and isolation from predators provided by caves and crevices seem essential. Others prefer the heat and dryness of sun-exposed roosts. Many species choose special nursery or hibernation roosts. Buildings are so widely exploited by bats (especially Vespertilionidae, Molossidae, and Emballonuridae) for regular diurnal or nursery roosts that the numbers of many species have probably become more abundant since the advent of architecture. Many bats also occupy nocturnal roosts, often rocky overhangs or cave entrances, for napping, for chewing food, or for shelter from bad weather.

Bats are usually colonial; indeed, some form very large

Allan Roberts



Indiana bats (*Myotis sodalis*) hibernating on a cave roof.

cave colonies. Generally, large colonies are formed by bats that roost in dense clusters, pressing against one another, although many roost spaced out, not touching. In trees, *Pteropus* may form outdoor camps numbering hundreds of thousands of individuals. Many species form smaller groups of several dozen to several hundred. Less commonly, bats are solitary; sometimes, the adult female roosts only with its most recent child. Occasionally one sex is colonial, and the other is apparently solitary. Some species regularly form mixed colonies (*e.g.*, *Mormoops* and *Chilonycteris* with *Leptonycteris*, *Monophyllus*, *Carollia*). The advantages of colonial or solitary life and the factors that govern colony size in bats with colonial predilection have not yet been established.

The roost requirements of many bats, rather precise in terms of light, temperature, and humidity, limit their distribution in space. Some of the Megachiroptera strikingly defoliate the trees on which they roost.

Elaborate communities of other animals are often satellites of cave-bat colonies. Among these are cave crickets, roaches, blood-sucking bugs, a variety of parasites (*e.g.*,

Speed in flight

fleas, lice, ticks, mites, and certain flies), and dermestid beetles and other insects that feed on cave-floor debris—guano, bat and insect corpses, and discarded pieces of food or seeds. Molds and other fungi are also conspicuous members of the cave-floor community. Bats and their excretions alter the cave environment by producing heat, carbon dioxide, and ammonia.

Migration. Many bats of temperate climates migrate annually to and from summer roosts and winter hibernation sites, with an individual often occupying the same roosts in seasonal sequence each year. Members of the same species may converge on a single hibernation cave or nursery roost from many directions, indicating that the choice of migration direction to and from these caves cannot be genetically determined. Migration time probably is genetically determined (*i.e.*, instinctive) and influenced also by weather conditions and the availability of food. Nothing is known of how migration goals are recognized or how their location is learned by succeeding generations. Female young, of course, are born at a nursery roost and may memorize its location, but how they know where to go at other times of year is not clear.

Female Mexican free-tailed bats migrate from Central Mexico to Texas and adjacent states each spring, returning south in the fall. Mating probably occurs in transient roosts during the spring flight. The migration is believed to remove pregnant (and lactating) females to a region of high food supply where they need not compete with males of their own species. Presumably they return to Mexico for its suitable winter climate and food supply and to meet their mates.

The North American red and hoary bats (*Lasiurus borealis* and *L. cinereus*) and the silver-haired bat (*Lasiurus noctivagans*) migrate in the fall from the northern U.S. and Canada to the southern states, returning in spring. Little is yet known of energy storage, navigation, or other specializations for migrations.

Orientation. Bats of the suborder Microchiroptera orient acoustically by echolocation ("sonar"). They emit short, high-frequency pulses of sound (usually well above the range of human hearing) and listen to the echoes returning from objects in the vicinity. By interpreting returning echoes, bats may identify the direction, distance, velocity, and some aspects of the size or nature or both of objects that draw their attention. Echolocation is used to locate and track flying and terrestrial prey, to avoid obstacles, and possibly to regulate altitude; orientation pulses may also serve as communication signals between bats of the same species. Bats of the megachiropteran genus, *Rousettus*, have independently evolved a parallel echolocation system for obstacle avoidance alone. Echolocation pulses are produced by vibrating membranes in the larynx and emitted via the nose or the mouth, depending upon species. Nose leaves in some species may serve to channel the sound.

The echolocation signals spread in three dimensions on emission, the bulk of the energy in the hemisphere in front of the bat or in a cone-shaped region from the nostrils or mouth. When the sound impinges on an intervening surface (an insect or a leaf, for example) some of the energy in the signal is reflected or scattered, some absorbed, and some transmitted and reradiated on the far side; the proportion of sound energy in each category is a function of wavelength and of the dimensions, characteristics, and orientation of the object. The reflected sound spreads in three dimensions, and some portion of it may impinge on the bat's ears at perceptible energy levels.

Bats' external ears are generally large, probably enhancing their value for detecting direction of incoming signals, and their middle and inner ears are specialized for high-frequency sensitivity. In addition, the bony otic (auditory) complex is often isolated acoustically from the skull, probably enhancing signal comparison by both ears. The threshold and range of hearing have been studied in several genera of bats, and, in each case, the region of maximum sensitivity coincides with the prominent frequencies of the outgoing echolocation signals.

The characteristics of echolocation pulses vary with family and even with species. Echolocation pulses of a sub-

stantial number of bat species have been analyzed in terms of frequency, frequency pattern, duration, repetition rate, intensity, and directionality. The prominent frequency or frequencies range from 12 kilohertz (one kilohertz is equivalent to 1,000 cycles per second) to about 150 kilohertz or more. Factors influencing frequency may include bat size, prey size, the energetics of sound production, inefficiency of propagation of high frequencies, and ambient noise levels.

Orientation pulses may be of several types. The individual pulse may include a frequency drop from beginning to end (frequency modulation, FM) or the frequency may be held constant (CF) during part of the pulse, followed by a brief FM sweep; either FM or CF pulses may have high harmonic content. The pulse duration varies with the species and the situation. In cruising flight, the pulses of the Asian false vampire (*Megaderma lyra*) are 1.5 milliseconds (0.0015 second), those of Dobson's mustache bat (*Chilonycteris psilotis*) 4 milliseconds, and those of the greater horseshoe bat 55–65 milliseconds. In goal-oriented flight, such as the pursuit of an insect or the evaluation of an obstacle or a landing perch, the pulse duration is systematically altered (usually shortened) with target distance, sometimes ending with pulses as short as 0.25 millisecond.

During insect pursuit, obstacle avoidance, and landing manoeuvres, there are three phases of pulse output design: search, approach, and terminal. The search phase, during which many bats emit about 10 pulses per second, precedes specific attention to a target. In the approach phase, which starts when the bat detects an object to which it subsequently devotes its attention, the bat raises the pulse rate to about 25 to 50 per second, shortens the pulses with decreasing distance, and often alters the frequency pattern. The terminal phase, which often lasts about 100 milliseconds, is characterized by extremely short pulses, repeated as rapidly as 200 or more per second, and ceases as the bat intercepts the target or passes it (the stimulus being, perhaps, the cessation of echoes); another search phase follows. During the brief terminal phase (a fraction of a second), the bat is engaged in final interception (or avoidance) manoeuvres and appears to pay little attention to other objects.

The use of echolocation to gain sensory information requires integration of the vocal and auditory centres of the brain, in addition to sensitive ears. Not only must the nervous system of the bat analyze in a few thousandths of a second the reflected, and thus altered, form of its own pulse, but it must separate this echo from those of other individuals and from others of its own pulses. All of this must be done while the animal (and often the target) is moving in space. In the laboratory, bats have been found to be able to identify, pursue, and capture as many as two fruitflies (*Drosophila*, about three millimetres [0.12 inch] long) per second, and to locate and avoid wires as fine as 0.1 or even 0.08 millimetre (0.004 or 0.003 inch) in diameter.

Research has provided some information on the mechanisms of bat sonar. There is evidence that the multiple frequencies of FM or harmonic patterns serve in determining target direction. The relative intensities of the various frequencies will be different at the two ears, allowing the animal to determine target direction when three or more frequencies are received. Target velocity may be measured by constant-frequency bats through the use of the Doppler shift, a change in perceived frequency due to the relative motion of the bat and its target. Changes in pulse-echo timing may provide information on target distance and velocity. The ratio of useful signal to background noise is increased by several mechanisms, including specializations of the middle ear and its ossicles (tiny bones), isolation of the cochlea (the area where sound energy is converted into nerve impulses), and adaptations of the central nervous system.

Food habits. Most bats feed on flying insects. In some cases, prey species have been identified from stomach contents or from discarded pieces under night roosts, but such studies have not yet provided an adequate measure of the spectrum of bat diets.

Separation
of male
and female
migrants

The
principle
of echo-
location

Capa-
bilities
of bat
sonar

Insects are identified and tracked in flight by echolocation. Large insects may be intercepted with the wing membranes and pulled into the mouth. Some moths, however, (Noctuidae, Arctiidae) are able to avoid capture by bats.

Some genera of bats (*Macrotus*, *Antrozous*, *Plecotus*, *Nycterus*) feed on arthropods, such as large insects, spiders, and scorpions, which they find on the ground, on walls, or on vegetation. These bats may either land on and kill their prey before taking off with it or pick it up with their teeth while hovering.

Three genera (*Noctilio*, *Pizonyx*, and *Myotis*) include at least one species that catches small fish and possibly crustaceans. All fish-eating species also feed on flying insects or have close relatives who do so. Each is specialized in having exceptionally large hindfeet armed with long, strong claws with which the fish are gaffed.

The Megachiroptera and many of the phyllostomatid genera feed on a variety of fruits, often green or brown in colour; usually such fruits are either borne directly on wood or hang well away from the bulk of the tree and have a sour or musky odour.

The pteropodid subfamily Eonycterinae (and some other fruit bats) and the phyllostomatid Glossophaginae feed, at least in part, on nectar and pollen. Many tropical flowers, adapted for pollination by these bats, open at night, are white or inconspicuous, have a sour, rancid, or mammalian odour, and are borne on wood, on pendulous branches, or held beyond or above the bulk of the plant. The phyllostomatid Phyllostominae may also feed on flowers.

Several phyllostomatid and megadermatid genera are carnivorous, feeding on small rodents, shrews, bats, sleeping birds, tree frogs, and lizards. The true vampires, which feed on the blood of large mammals or birds, land near a quiet prospective victim, walk or jump to a vulnerable spot on it where the skin is relatively exposed—the edge of the ear or nostril, around the anus, or between the toes, for example—make a scooping, superficial bite from which the blood oozes freely, and lap the blood with very specialized tongue movements. Each vampire requires about 15 millilitres (about one cubic inch) of blood per night.

The interaction of bats with their food, be it insects, fruit, or flowers, probably has a substantial impact on some biological communities. Many plants are dependent on bats for pollination; other plants benefit from seed dispersal by bats. Moths of two families are known to take evasive or protective action on hearing bat pulses nearby, an adaptation that implies heavy predation.

Maintenance behaviour. Bats are meticulous in their grooming, spending a fair part of the day and night combing and grooming their fur and cleansing their wing membranes. Generally they comb with the claws of one foot, while hanging by the other; they remove the combings and moisten their claws with their lips and tongue. On the wing membranes, in particular, they use the mouth meticulously, perhaps oiling the skin with the secretions of dermal (skin) glands while cleansing it.

Social interactions. Bats often segregate by sex. As noted, in many species, pregnant females occupy special nursery roosts until their young are independent. In some species, the sexes occupy the same general roost but gather in separate clusters. In others, the sexes intermingle or arrange themselves into a pattern within a group—the females centrally, for example, and the males peripherally. Sexual segregation during foraging has been reported for several species. Among bats that migrate over long distances, such as Mexican free-tailed, red, and hoary bats, the sexes may meet only briefly each year.

FORM AND FUNCTION

Anatomical specializations. Bats are mammals with front limbs modified for flight. The fingers, other than the thumb, are greatly elongated and are joined by a membrane that also extends from the posterior border of the forearm and upper arm to the side of the body and leg as far as the ankle or foot. Only the thumb, and occasionally the index finger, end with a claw. The wing membrane consists of two layers of skin, generally darkly pigmented and naked, between which course blood vessels and nerves.

When not fully extended, the wing skin is gathered into wrinkled folds by elastic connective tissue and muscle fibres. Some of the fingers, especially the third, fold over when the bat is not in flight; the wing may then be quite tightly folded or may partly enfold the bat's undersurface. The thumb, always free of the wing membrane, is used for walking or climbing in some species; in others, it is used for handling food. Bats that walk often have pads or suction disks on their thumbs or wrists or both, and many female bats use the thumbs for suspending themselves, hammock fashion, when giving birth.

Wing shape, governed by the relative lengths of the forearm and the fingers, varies greatly, in adaptation to flight characteristics.

Most bats have a membrane, consisting of skin like

Drawn by R. Keane, based on *Natural History* (October 1958)

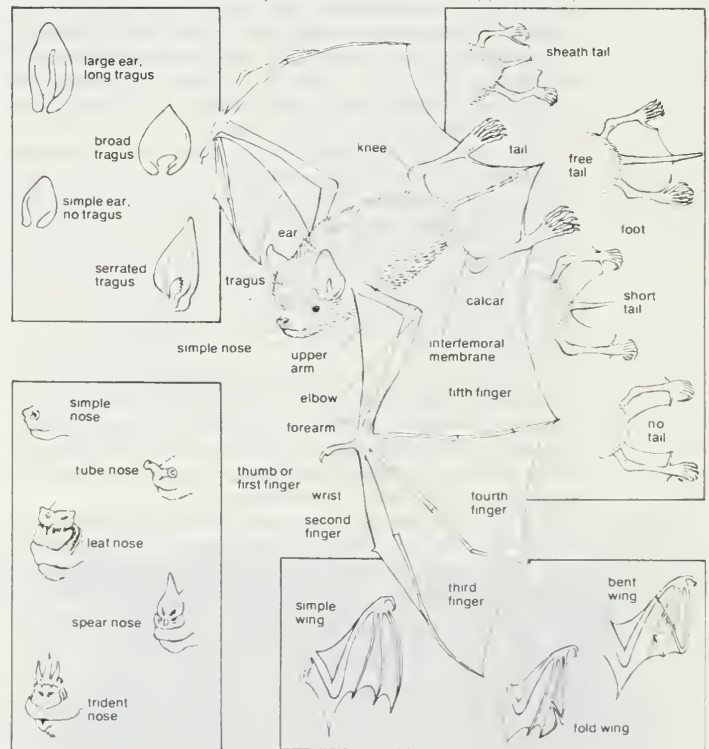


Figure 18: Typical microchiropteran bat (*Myotis*). Insets show variations of structures found in other bats.

the wings, that extends between their legs (intermembral membrane). In the midline, the intermembral membrane is usually supported, at least in part, by the tail, and the distal edges are often shaped in flight by greatly elongated heel bones. The intermembral membrane, especially well developed in insectivorous, carnivorous, and fish-eating bats, is less well developed or even absent in the vampires and in fruit- and flower-feeding bats. Many bats in flight, on catching large prey, bring the membrane forward and by flexing the neck and back tuck the prey against and into the membrane while taking a new tooth hold. By this manoeuvre, the bat takes hold of the victim head first and is able to kill or disable it promptly.

At rest, a bat's head is its most striking feature. The pinna (projecting portion) of the external ear usually is extremely large and often funnel shaped. In several genera that feed on terrestrial arthropods, the ears are particularly oversized, probably for ultraprecise directional assessment. The tragus (a projection on the anterior side of the auditory canal) or antitragus (on the posterior side) may also be conspicuous. The ears are often highly mobile, sometimes flicking back and forth in phase with the production of sonar signals. In some species the ears are immobile, but, in all cases, they probably function jointly for directional analysis.

Bats often have a muzzle of somewhat rodent-like or foxlike proportions but, in many, the face has a pushed in, pug-dog appearance. In the nectar feeders, the snout is elongated to house the long, extensible tongue. Many bats

Fishing by bats

Structure of the wing

Shape of the face

have a facial ornament, the nose leaf, consisting of skin and connective tissue surrounding the nostrils and extending as a free flap or flaps above the nostrils and in front of the face. The complexity and shape of the nose leaf varies with family; its presence correlates with nasal emission of orientation signals. Thus, it is supposed that the nose leaf influences sound output, perhaps by narrowing the beam, but evidence is sparse.

A bat's neck is likely to be short and relatively immobile, the chest and shoulders large and well muscled, the hips and legs slender.

Most bats are well furred (except for the wing membranes), generally in shades of brown, tan, gray, or black on top and in lighter shades ventrally. Red, yellow, or orange variants occur in many species. Among species that roost in the open, speckled or mottled fur, bright or light-coloured spots or stripes, or bright red, yellow, or orange shading on the pendant head, neck, and shoulders are common. Mottled fur may blend with lichen-covered bark or rock. Bright spots may simulate the speckled sunlight of the forest canopy as seen from below. Stripes probably break up contours. The colouring seen while the animal is hanging may be a kind of countershading for concealment or it may enhance the bat's simulation of a ripening fruit or a dead leaf. Many bats who roost externally hang from a branch by one foot, which then looks like a plant stem.

Many bats have large dermal (skin) glands (the location of which depends on family), which secrete odorless substances that may serve as species or sex recognition signals. Some glands may also supply oils for conditioning the skin or waterproofing the pelage.

Thermoregulation. Although some bats maintain fairly even body temperatures, a large number undergo periodic raising or lowering of their temperature. Many of the vespertilionids and rhinolophids and a few molossids drop their body temperature to ambient temperature shortly after coming to rest (a condition called heterothermy). They raise their temperature again on being aroused or when readying themselves for nocturnal foraging. When fully active, bats have a body temperature of about 37° C (98.6° F). During pregnancy, lactation, and juvenile growth, bats probably thermoregulate differently, more closely approximating stability. The drop in body temperature, if the ambient temperature is relatively low, results in the bat's assuming a lethargic state. Energy is conserved by "turning down the thermostat," but the bat is rendered relatively unresponsive to threats by predators or weather. Heterothermic bats, therefore, generally roost in secluded, isolated sites, often in crevices, and have specializations for arousal. One or more sensory systems and the brain remain sensitive at low temperatures and initiate the necessary heat production for arousal. Heat is generated by the metabolism of brown fat and by shivering.

Many bats that exhibit daily torpor hibernate during the winter and must, therefore, store energy as body fat in the fall, increasing their weight by 50 to 100 percent. They must also migrate from the summer roost to a suitable hibernation site (often a cave); *i.e.*, one that will remain cool and humid but will not drop below freezing throughout the winter. Large populations often aggregate in such caves. Hibernation also involves the suspension of temperature regulation for long periods; adaptations of circulation, respiration, and renal function; and suspension of at least some aspects of diurnal activity. Bats of hibernating species generally court and mate in the fall when the males are at their nutritional peak.

Bats of several tropical families are homeothermic. A spectrum of degrees of homeothermy and heterothermy probably will be discovered in the order.

Digestion and water conservation. Digestion in bats is unusually rapid. They chew and fragment their food exceptionally thoroughly and thus expose a large surface area of it to digestive action. They may begin to defecate 30 to 60 minutes after beginning to feed, thus reducing the load that must be carried in flight.

Some bats live in sun-baked roosts without access to water during the day. They may choose these roosts for their heat, thus conserving body heat, but it is not yet known how they hold their body temperature down without using

water. In the laboratory, bats die if body temperature rises above about 40–41° C (104°–106° F).

Senses. Bats have been considered in folklore to be blind. In fact, the eyes in the Microchiroptera are small and have not been well studied as yet. Among the Megachiroptera, the eyes are large, but vision has been studied in detail only in *Pteropus*. These bats are able to make visual discriminations at lower light levels than can man. The Megachiroptera fly at night, of course, and some genera fly below or in the jungle canopy where light levels are very low. Except for *Rousettus*, none are known to orient acoustically.

Studies of several genera of Microchiroptera have revealed that vision is of use in long distance navigation and that obstacles and motion can be detected visually. Bats also presumably use vision to distinguish day from night and to synchronize their internal clock with the local day-night cycle.

The senses of taste, smell, and touch in bats do not seem to be strikingly different from those of related mammals. Smell is probably used as an aid in locating fruit and flowers and possibly, in the case of vampires, large vertebrates. It may also be used for locating an occupied roost, colleagues of the same species, and individuals of the correct sex. Many bats depend upon touch, aided by well-developed facial and toe whiskers and possibly by the projecting tail, to place themselves in comforting body contact with rock surfaces or with colleagues in the roost.

EVOLUTION AND PALEONTOLOGY

The fossil record of bats prior to the Pleistocene (about 2,500,000 years ago) is limited and reveals little about bat evolution. Most fossils can be attributed to living families. Skulls and teeth compatible with early bats are known from the Paleocene (about 60,000,000 years ago), but these fossils may equally well have been from insectivores, from which bats are clearly separable only on the basis of adaptations for flight. By the middle Eocene (45,000,000 years ago) bats with full flight had evolved.

The Chiroptera are readily divided into two suborders—Megachiroptera and Microchiroptera. The Megachiroptera orient visually and exhibit a number of primitive skeletal features. The Microchiroptera orient acoustically. It is not certain that they have a common origin. The suborders either evolved separately from flightless Insectivora or diverged very early in chiropteran history.

The two principal geographic centres of bat evolution appear to be the Australo-Malaysian region, with about 290 species, and the New World tropics, with about 230 species. Comparable ecological niches in the Old and New World are largely occupied by different genera of bats, usually of different families.

CLASSIFICATION

Distinguishing taxonomic features. The order Chiroptera is defined by true flight. The elongated finger bones support the wing membrane. Bats are also characterized by their generally small size, marked pectoral specialization for flight, and relatively weak pelvic and leg development. The ulna is reduced, claws are absent on fingers except for the thumb (and occasionally the second finger), and the knee is directed posteriorly and laterally. The maximum complement of milk teeth is 22 and of permanent teeth, 38; the minimum of permanent teeth is 20. The dental formula indicates the number of pairs of upper and lower incisors, canines, premolars, and molars, respectively, and total number of teeth.

Annotated classification.

ORDER CHIROPTERA

Eighteen living families, 174 genera, and about 900 living species.

Suborder Megachiroptera

One family.

Family Pteropodidae (flying foxes and Old World fruit bats)

Generally large, fruit- or flower-feeding; lack acoustic orientation (except *Rousettus*); ears small, eyes large, vision well-developed; generally roost in trees, often colonial; often show countershading, cryptic markings, or bright fur colours or patterns. Index finger generally clawed, tail short or lacking, inter-

Use of
the eyes

Lowering
of body
temper-
ature

Two main
evolu-
tionary
lines of
bats

femoral membrane reduced. Muzzle simple in appearance (except *Hypsipnathus*). Generally cannot walk but can move along branches in hanging posture. Forearm length varies from 37 mm (1.46 inches; *Macroglossus*) to 220 mm (8.66 inches; *Pteropus vampyrus*). Teeth modified for fruit- and flower-feeding. Dental formula $\frac{(1-2) \cdot 1 \cdot 3 \cdot (1-2)}{(0-2) \cdot 1 \cdot 3 \cdot (2-3)} = 24-34$. Old World tropics and subtropics, including many Pacific islands; 39 Recent genera, about 154 species.

Suborder Microchiroptera
Seventeen families.

Family Rhinopomatidae (mouse-tailed bats)

Small, insectivorous. Tail very long and largely free beyond a narrow interfemoral membrane, forearm very long, ears large, small nose leaf, primitive shoulder girdle. Dental formula $\frac{1 \cdot 1 \cdot 1 \cdot 3}{2 \cdot 1 \cdot 2 \cdot 3} = 28$. Store fat (probably on a seasonal basis). Roost on vertical surfaces, probably not in total darkness or isolation. Tropical distribution from northern Africa through southern Asia as far as Sumatra; 1 genus, 3 species. Generally considered to be the most primitive of all living Microchiroptera.

Family Emballonuridae (sheath-tailed or sac-winged bats)

Small to medium size. Ears large but simply shaped, eyes small, muzzle sharp but plain; tail short, perforating dorsal surface of well-developed interfemoral membrane. Several genera have a glandular pouch in the wing extension, anterior to the arm. Relatively unspecialized shoulder girdle and arm articulation. Dental formula $\frac{2 \cdot 1 \cdot 2 \cdot 3}{3 \cdot 1 \cdot 2 \cdot 3} = 34$, $\frac{1 \cdot 1 \cdot 2 \cdot 3}{3 \cdot 1 \cdot 2 \cdot 3} = 32$, or $\frac{1 \cdot 1 \cdot 2 \cdot 3}{2 \cdot 1 \cdot 2 \cdot 3} = 30$. Insectivorous. Roost on vertical surfaces, such as tree trunks, cliff faces, cave entrances, and walls; some favour buildings, especially belfries. Some densely colonial but not touching one another; others form small groups or are solitary. Hang suspended from toes with wrists propped against wall. Worldwide tropical distribution (excluding West Indies and some Pacific islands); 12 genera, about 50 species; each genus restricted to either Old or New World.

Family Nycteridae (slit-faced or hollow-faced bats)

Small to medium size. The humerus and pectoral girdle not greatly specialized, skull with peculiar nasal fossa (depression), cleft nose leaf, and a deep midline facial cleft behind and above the nostrils. Ears large, wings broad, tail long with bifid (split) end, calcars (heel bones) greatly elongated, tail and calcars supporting well-developed interfemoral membrane. Dental formula $\frac{2 \cdot 1 \cdot 1 \cdot 3}{3 \cdot 1 \cdot 2 \cdot 3} = 32$. Insectivorous, mostly preying on terrestrial forms or those resting on vegetation, rocks, or walls. Cannot walk. Roosts usually dark and humid, some species roosting externally in jungle canopy. Generally form small nontouching colonies, but some are solitary. Distributed through most of tropical Africa, Malaysia, and Indonesia; 1 genus, 13 species.

Family Megadermatidae (false vampires)

Moderately large bats. External ears very large and fused across midline; tragus bifid; nose leaf large with truncated end; eyes relatively large; wings broad, interfemoral membrane well-developed and supported distally by heel bones, no external tail. Females bear false inguinal nipples. Premaxillae lacking; dental formula $\frac{0 \cdot 1 \cdot (1-2) \cdot 3}{2 \cdot 1 \cdot 2 \cdot 3} = 26-28$. Insectivorous, principally on terrestrial arthropods, as in Nycteridae; at least 2 species, *Megaderma lyra* and *Macroderma gigas*, also feed on small vertebrates hunted and taken in the same fashion as arthropod prey. Cannot walk. Form small nontouching colonies usually in dark, secluded caves or abandoned buildings. *Lavia frons* are at least partly diurnal and roost in trees in the savanna and open forest. Central Africa, Southeast Asia, and Australia; 4 genera, 5 species.

Family Hipposideridae (Old World leaf-nosed bats)

Small to large bats. Complex nose leaf with subordinate leaflets and compartments; large ears, widely separated and highly mobile, antitragus well-developed; interfemoral membrane well-developed, tail generally projecting slightly beyond distal edge of membrane. Often bear glandular pouch on forehead; females have false inguinal nipples. Dental formula $\frac{1 \cdot 1 \cdot 2 \cdot 3}{2 \cdot 1 \cdot 2 \cdot 3} = 30$. Colour usually drab brownish but red phases not uncommon. Insectivorous, usually on flying insects. Mostly colonial, nontouching, and roosting in humid caves, tree hollows, culverts, or buildings. A few roost externally in branches of trees, a few are solitary. Do not walk. Old World tropics; 9 genera, about 60 species, the genus *Hipposideros* particularly successful in numbers of species and individuals throughout the Old World tropics.

Family Rhinolophidae (horseshoe bats)

Small to moderately large size. Complex nose leaf; large, highly mobile ears, well-developed antitragus; wings short and rounded; well-developed interfemoral membrane, supported by tail; calcanea (backs of heels) weak. Fur generally brown (red phases occur). Females bear false inguinal nipples. Dental formula $\frac{1 \cdot 1 \cdot 2 \cdot 3}{2 \cdot 1 \cdot 3 \cdot 3} = 32$. Dark, humid roosts selected, especially caves, but tree hollows, buildings, and culverts as well. Generally colonial, nontouching. Cannot walk. Insectivorous, usually on flying insects. Old World, including parts of western Europe, Central Asia, and Japan; 2 genera, about 70 species, the genus *Rhinolophus* one of the most successful in species and numbers.

Family Noctilionidae (bulldog bats)

Medium size. Muzzle heavy but unadorned; lips full; internal cheek pouches; ears large, pointed, and mobile; wings long and narrow. Dental formula $\frac{2 \cdot 1 \cdot 1 \cdot 3}{1 \cdot 1 \cdot 2 \cdot 3} = 28$. Tail well-developed, extending to midpoint of large interfemoral membrane, which is pierced dorsally by tail tip, membrane supported distally by very well-developed calcars and calcanea. Feet large, or very large (*N. leporinus*). Colour dark brown to rufous orange dorsally. Musky odour. Walk well, often roost in crevices, tree hollows, attics, grottoes, and caves; colonial, in touching clusters. *Noctilio labialis* feeds on flying insects, as does *N. leporinus*, which also gaffs fish. Tropical America; 1 genus, 2 species.

Family Mormoopidae

Small. Insectivorous on flying insects. Some walk. All lack nose leaf but have elaborate lip leaves. Colour ranges from brown through orange, red, and yellow. Densely colonial in dark caves, colonies often numbering tens of thousands. Tail and interfemoral membrane well-developed. Dental formula $\frac{2 \cdot 1 \cdot 2 \cdot 3}{2 \cdot 1 \cdot 3 \cdot 3} = 34$. Tropical Central and South America; 3 genera, 9 species.

Family Phyllostomatidae (American leaf-nosed bats)

Small to large size. Nose leaf simply shaped, ears often large and mobile, wings generally short and broad, tail and interfemoral membrane quite varied (from absent to well-developed); dental formula varied, from 26 to 34; fur colour and patterns varied. Insectivorous (the insects eaten may be flying or terrestrial types), carnivorous, fruit- and flower-feeding species. Generally do not walk. Colonial, often densely so; generally roosting in touching clusters in caves, tree hollows, buildings, or culverts or in the open under bridges or eaves, in the crests of palm trees, or on the underside of palm leaves. Flight swift and straight to hovering. Southwestern U.S. through tropical America; 47 genera, about 120 species.

Subfamily Phyllostomatinae. Medium to large. Nose leaf often of striking size, ears large. Insectivorous, fruit-feeding, or carnivorous. Varied dental formulas. Subfamily may be an artificial grouping.

Subfamily Glossophaginae. Small to medium size. Small nose leaf, snout elongated; tongue elongated and extensible. Wings broad; hovering flight. Feed on pollen, nectar, fruit.

Subfamilies Carollinae, Strumirinae, Stenoderminae. Medium size. Fruit-feeding; some may be insectivorous. Nose leaves well-developed. Many Stenoderminae have white or light facial stripes, may roost in the open; 2 species alter palm leaves as roosts.

Subfamily Phyllonycterinae. Small bats endemic to the West Indies, probably fruit- and flower-feeding. Legs and feet exceptionally well-developed and body flexible. Beige, maize, or light brown. Nose leaves very small; snout and tongue moderately elongated. Cave-dwelling.

Family Desmodontidae (vampire bats)

Medium-sized bats. Small nose leaf. Teeth highly specialized for cutting skin, cheek teeth reduced; dental formula in *Desmodus* $\frac{1 \cdot 1 \cdot 1 \cdot 1}{2 \cdot 1 \cdot 2 \cdot 1} = 20$. Hindlegs and thumbs very well adapted for walking and jumping. Tail absent; interfemoral membrane reduced. Feed on blood of large mammals or birds. Roost in caves, hollow trees, and culverts; colonial. Most of tropical America, excluding West Indies; 3 genera, 3 species.

Family Natalidae (funnel-eared bats)

Small, slenderly built. Gray, buffy, yellow, or reddish; fur deep. Well-developed tail and interfemoral membrane. Ears large; snout plain; dental formula $\frac{2 \cdot 1 \cdot 3 \cdot 3}{3 \cdot 1 \cdot 3 \cdot 3} = 38$. Natalids walk clumsily but do not enter crevices; they are cave-dwelling, colonial in nontouching groups. Feed on flying insects. Central America, and northern South America, West Indies; 1 genus, 4 species.

Family Furipteridae (smoky bats)

Small, delicately built. Thumb vestigial; snout plain; tail long, ending short of distal edge of well-developed interfemoral membrane; legs long; feet small. Dental formula $\frac{2 \cdot 1 \cdot 2 \cdot 3}{3 \cdot 1 \cdot 3 \cdot 3} = 36$. Biology unknown; probably insectivorous. Northern South America; 2 genera, 2 species.

Family Thyropteridae (disk-wing bats)

Second finger reduced to rudiment, base of thumb and sole provided with sucking disk, simple muzzle, ears large; dental formula $\frac{2 \cdot 1 \cdot 3 \cdot 3}{3 \cdot 1 \cdot 3 \cdot 3} = 38$. Insectivorous, roost alone or in small groups, often in still furled banana leaves. Biology poorly known. Central America and northern South America, excluding West Indies; 1 genus, 2 species.

Family Myzopodidae (Old World sucker-footed bats)

Small, plain muzzle, large ears with peculiar mushroom-shaped lobe, dental formula $\frac{2 \cdot 1 \cdot 3 \cdot 3}{3 \cdot 1 \cdot 3 \cdot 3} = 38$. Thumb and sole with adhesive disks, vestigial thumb claw; tail extends free beyond interfemoral membrane; specialized scapulo-humeral articulations. Probably insectivorous (biology unknown); endemic to Madagascar; 1 species.

Family Vespertilionidae (common bats)

Small to medium sized. Muzzle plain; eyes small; ears moderate to large, tragus well-developed; dental formula varied. $\frac{1 \cdot 1 \cdot 1 \cdot 3}{2 \cdot 1 \cdot 2 \cdot 3} = 28$ to $\frac{2 \cdot 1 \cdot 3 \cdot 3}{3 \cdot 1 \cdot 3 \cdot 3} = 38$. Wings generally long and moderately narrow; tail and interfemoral membrane well-developed. All walk well, often entering crevices. Insectivorous, some on flying, some on terrestrial insects; a few genera (*Piponyx* and some *Myotis*) eat fishes. Mostly roost in caves, attics, barns, hollow trees, boulder heaps, the twig work of birds' nests, or roof thatching; a few (*Lasiurus*) roost in branches, on tree trunks, or in hollow core of bamboo (*Tylonycteris*). Generally colonial in dense, touching clusters; a few solitary. Many temperate species hibernate and migrate; many drably coloured; several that roost externally are either spotted (*Enderma*) or speckled (*Lasiurus*). Family worldwide to tree line, including many oceanic islands; some genera (*Myotis*, *Eptesicus*, and *Pipistrellus*) also worldwide; 35 genera, about 290 species.

Family Mystacinidae (New Zealand short-tailed bats)

Small, with simple vespertilionid-like head; highly adapted for walking. Wings fold exceptionally compactly; thumb and toe claws long and sharp; tail perforates interfemoral membrane dorsally. Dental formula $\frac{1 \cdot 1 \cdot 2 \cdot 3}{1 \cdot 1 \cdot 2 \cdot 3} = 28$. Feeds on flying and terrestrial insects. Biology poorly known. New Zealand; 1 species.

Family Molossidae (free-tailed bats)

Robustly built, small to very large. Tail projects well beyond well-developed interfemoral membrane; ears large, rather immobile, often fused to one another, and of very unusual shapes; lips and snout often heavy, eyes tiny. Wings very long and narrow; legs well-developed for walking; toes often bear bristles; dental formula varied from $\frac{1 \cdot 1 \cdot 1 \cdot 3}{1 \cdot 1 \cdot 2 \cdot 3} = 26$ to $\frac{1 \cdot 1 \cdot 2 \cdot 3}{3 \cdot 1 \cdot 2 \cdot 3} = 32$. Often have a conspicuous odour; fur short, usually dark brown or black. Many highly colonial, as many as millions clustering in dense, touching groups; roost pressed against fellows or walls, often in crevices. Occupy caves, tree hollows, and buildings. Walk exceptionally well; many prefer hot, dry roosts. Feed on flying insects; some migrate. Worldwide in tropics and subtropics, with a few species ranging into mild temperate regions; 11 genera, about 90 species (the genus *Tadarida* worldwide).

(A.N.)

Primates (lemurs, lorises, tarsiers, monkeys, apes, and hominids, including man)

The mammalian order Primates, which includes the prosimians, monkeys, apes, and man, has long excited human interest, because it contains man's closest relatives. The obvious similarity of monkeys and apes to man was noted long before any genetic relationship was postulated. Recognizing that monkeys possess intelligence substantially exceeding that of other familiar mammals, man has valued these animals as companions and pets for centuries. In recent decades, in addition to becoming increasingly popular as pets, nonhuman primates have served as substitutes for man in situations, such as experimental medicine and space science, which require near-human biological reactions.

The two main groups of the order Primates are the prosimians (lemurs, lorises, and tarsiers) and the anthropoids (monkeys, apes, and man). As classified at present, the suborder Prosimii consists of six families; the Tupaiidae (tree shrews, though these are considered by some authorities to belong in the order Insectivora, a view elaborated in the section *Insectivora*). Lemuridae (lemurs), Indriidae (indri, sifaka, and avahi), Daubentonidae (aye-aye), Lorisidae (galagos and lorises), and Tarsiidae (tarsiers). Collectively, the prosimian suborder is frequently referred to as lower primates.

The suborder Anthropoidea, sometimes called "higher" primates, also comprises six families: Callitrichidae (marmosets and tamarins), Cebidae (South American monkeys other than marmosets), Cercopithecidae (African and Asian monkeys), Hylobatidae (lesser apes; siamangs and



Figure 19: Representative prosimians.

gibbons), Pongidae (great apes: orangutan [orangutan], gorilla, and chimpanzee), and Hominidae (men, living and extinct). The two geographically separated stocks, the Old World and the New World monkeys, are often referred to as catarrhines and platyrrhines, respectively, terms which derive from the shape of the nose (see below, *Form and function*).

GENERAL CONSIDERATIONS

Historical background of primate studies. The inclusion of man in the order provides the rationale for the vigour with which this group has been studied by scientists since the time of Galen of Pergamum. Aristotle and Hippocrates, in the 3rd and 4th centuries BC, recognized the similarity of man and apes, but it was Galen who demonstrated the truth of this kinship by dissection. He wrote, "the ape is likeliest to man in viscera, muscles, arteries, veins, nerves and in the form of bones." It should be noted that Galen was in fact referring to monkeys, not the true apes, which were unknown to western man until the 15th century. None of these early scientists saw any evolutionary significance in the similarity of man and "apes," a correspondence that they regarded as purely coincidental. An inkling of the truth of man's relationship with primates must have penetrated the mind of St. Albertus Magnus, probably the best leading naturalist of the Middle Ages, who produced a classification of animal life in his book *De animalibus*. Albertus' classification, which placed man between "apes" (monkeys) on the one hand and "animals" on the other, provides the first whiff of the "missing-link" concept, which later was to befog the issue of man's place in nature.

The Dark Ages were aptly named as far as knowledge of primates is concerned. The first evidence of a renaissance of interest was in the time of Vesalius, the great Belgian anatomist of the 16th century, who published a comparative anatomy of man and "apes" in order to confound the precepts of Galen. He did not succeed in disproving Galen's assertion that "apes were likeliest to man" but, unwittingly, he succeeded in stirring up an interest in the biology of primates that has never since flagged. The first true ape studied as a scientific specimen was a chimpanzee dissected by Edward Tyson, an English anatomist, in 1699. Tyson's specimen, which he called the "Orang-Outang, sive Homo Sylvestris," is to this day housed in the British Museum (Natural History), mounted in a standing position that reflects Tyson's belief that he had discovered the pygmy, a race of humans known since the time of the ancient Greeks. Tyson wrote of his "pygmie" that it was "no man, nor yet a common ape but a sort of animal between both." It never occurred to Tyson or his contemporaries, who believed that all animals had been created independently in their current image, that man, apes, and monkeys were connected by common evolutionary descent. In 1758, Linnaeus—the father of animal and plant classification—added the lemurs and bats to the monkeys, apes, and man and called the whole assemblage the Primates. It is to be noted, however, that Linnaeus was sufficiently perceptive to see that man was a primate. His conclusion was regarded as a grave blow to human dignity, and it was followed by new classifications such as that of Blumenbach in 1776, placing man in a separate order. Man was not again considered part of the primate order until a century later when the English anatomist St. George Mivart in the climate of post-Darwinian thought published his classification of primates.

The first evolutionist was a French scholar of the late 18th century, the Chevalier de Lamarck, who saw animal life as an uninterrupted continuity in which old species were transformed into new species in a sequence of increasing complexity and perfection. In 1821 Baron Georges Cuvier, a rabid anti-evolutionist, had the historic distinction of describing *Adapis*, the first fossil primate genus ever recognized. Fossils such as *Adapis*, Cuvier believed, were the remains of animals destroyed by past catastrophes (floods, earthquakes, etc.) and living animals were new stocks divinely created to fill the vacuum, a view consistent with

the widely held notion that species were immutable. During the early 19th century a number of geologists and biologists questioned the doctrine of immutability, but it was not until 1859, with the publication of Charles Darwin's "On the Origin of Species by Means of Natural Selection," that positive evidence was provided, along with a sound alternative theory. The Darwinian contention that not only had man evolved, but he had evolved from a simian ancestor resulted in acrimonious debate among scientists, theologians, philosophers, and laymen. As influential zoologists and anatomists rose to support Darwin, the truth of man's primate consanguinity began to be accepted, if not actually relished. Today, few scientists deny that man and the lower primates belong in the same order; in fact, much current research is directed toward closing the apparent gap between the highest of the nonhuman primates, the chimpanzee and gorilla, and man.

In times past, the public image of primates was largely dictated by prevailing religious beliefs. In Asian countries, where primates abound, monkeys have for a long time been regarded with various degrees of deference that, among Hindus in India, for instance, amounts to worship. In Europe and North America, where monkeys and apes are totally absent, no religious sect has attached divine significance to them and, in fact, the reverse has been the case, monkeys at various times having been regarded as the personification of evil and depravity, familiars of the devil. This image, however, is fading as a result of enlightened instruction in schools and the advances in naturalistic presentation of primates in zoos. The nonhuman primate is becoming generally accepted by laymen as an animal of peculiar interest to man with many amusing and endearing qualities.

The qualities of primates, however, are less endearing to farmers and agriculturists in certain parts of the world. In South Africa, the chacma baboon (*Papio ursinus*) competes with domestic sheep for grazing lands and is an occasional predator of lambs; in West and Central Africa native crops are subject to daily assaults by forest-living monkeys; and, in India, macaques, which have been accorded a semi-sacred status, live alongside man in towns and villages and are parasitic upon him for their food and shelter.

Scientific interest in nonhuman primates, their structure at all levels and their way of life, is currently in the ascendancy. Their value as research animals increases year by year, so that at present over a quarter of a million wild monkeys are consumed annually by laboratories in the study of human diseases, in the production of vaccines, in the experiments of organ transplantation, in the testing of drugs and, even, for the clinical trials of new cosmetics. Clearly their scientific usefulness raises important problems of conservation of primate stocks in the wild. Other disciplines whose researches depend upon observations and upon experimentation with nonhuman primates include those of endocrinology, neurology, psychology, and sociology. Much is being learned as a result of such studies that is of great significance for man and his future betterment.

Size range and adaptive diversity. Members of the order Primates show a remarkable range of size and adaptive diversity. The smallest primate (in linear dimensions) is probably the dwarf or mouse lemur (*Microcebus murinus*) of Madagascar, which measures approximately 13 centimetres (about five inches) in head and body length; the most massive is certainly the gorilla (*Gorilla gorilla*), whose head and torso length may reach 100 centimetres (39 inches; the comparable dimension in man may exceed this figure, but the gorilla is considerably more robust). In terms of body weight, the range is even more striking. The pen-tailed tree shrew (*Ptilocercus lowii*) weighs on average 46 grams (1.62 ounces); the gorilla's weight may be more than 3,000 times as great, varying from 140 to 180 kilograms (about 300–400 pounds).

In terms of habitat, primates are tropical animals and occupy two major vegetational zones: the tropical forest and woodland-grassland vegetational complexes. Each of these zones has produced in its resident primates the appropriate adaptations, which are not remarkable for their

Discovery of the first true ape

Religious views toward primates

Size and weight

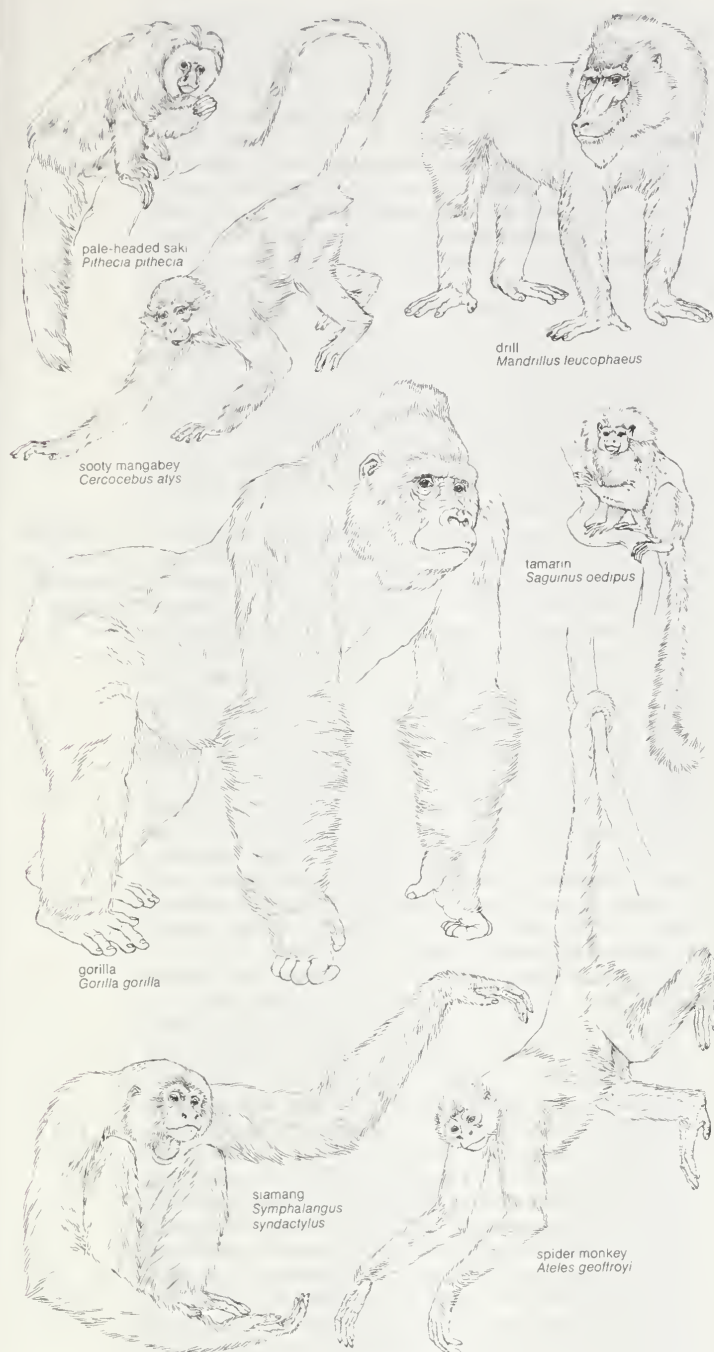


Figure 20: Body plans of representative anthropoids.

Drawing by R. Keane

diversity. There is perhaps more diversity of bodily form to be seen amongst forest-living species than among savanna inhabitants. One of the explanations of this difference is that it is the precise pattern of locomotion rather than the simple matter of habitat that governs overt bodily adaptations. Within the forest there are a number of ways of moving about. An animal can live on the forest floor or in the canopy, for instance, and within the canopy it can move in three particular ways: (a) by leaping—a function principally dictated by the hindlimbs; (b) by arm swinging (brachiation)—a function particularly of the forelimbs; (c) by quadrupedalism—a function equally divided between the forelimbs and the hindlimbs. On the savanna, or in the woodland-savanna biome, which substantially demands adaptations for ground-living locomotion rather than those of tree-living, the possibilities are limited. If bipedal man is discounted, there is a single pattern of ground-living locomotion, which is called quadrupedalism. Within this category there are at least two variations on the theme: (a) knuckle-walking quadrupedalism, and (b) digitigrade quadrupedalism. The former gait is characteristic of the African apes (chimpanzee and gorilla), and the latter of

baboons and macaques, who walk on the flats of their fingers. After man, the Cercopitheciinae are the most successful colonizers of nonarboreal habitats.

The structural adaptations of primates resulting from locomotor differences are considered below in more detail, but they do not prove to be very extensive. Primates are a homogeneous group in terms of morphology, and it is only in the realm of behaviour that differences between primate taxa are clearly discriminant. It can be said that the most successful primates (judged in terms of the usual criteria of population numbers and territorial spread) are those that have departed least from the ancestral pattern of structure, but furthest from the ancestral pattern of behaviour. "Manners maketh man" is true in the widest sense of the word; in the same sense, manners delineate primate species.

Distribution and abundance. The nonhuman primates have a wide distribution throughout the tropical latitudes of Africa, India, Southeast Asia, and South America. Within this tropical belt, which lies between latitudes 25° N and 30° S, they have a considerable altitudinal range. In Ethiopia the gelada baboon (*Theropithecus*) is found living at altitudes up to 5,000 metres (16,000 feet). Mountain gorillas of the Virunga Volcanoes are known to travel across mountain passes at altitudes of more than 4,200 metres (13,800 feet) when travelling from one high valley to another. The howler monkeys of Venezuela (*Alouatta seniculus*) live at 2,500 metres (8,000 feet) in the Cordillera da Merida and in northern Colombia the dourocouli (*Aotus*) is found in the tropical montane forests of the Cordillera Central.

A few primate species occur naturally in temperate latitudes. The rhesus monkey (*Macaca mulatta*), for example, is found near Peking, and Japanese macaques (*Macaca fuscata*) are to be found on the Shimokita Peninsula, a hatchet-shaped promontory of the Japanese island of Honshu on latitude 41° N, where snow covers the ground for four or five months of the year. Macaques are particularly hardy animals and can withstand cold climates in the wild and captivity extremely well. The limiting factor to the occupation of temperate zones by monkeys is the winter food supply, combined with the short daylight hours available for foraging.

Primates in general are still plentiful in the wild; for the majority of common species there are few serious conservation problems at present. A number of species, particularly the orang-utan, the gorilla, some of the Madagascan lemurs, and some South American species, such as the golden marmoset, however, are in serious danger of extinction unless their habitats can be preserved in perpetuity and predatory habits of man kept under control.

NATURAL HISTORY

Reproduction and life cycle. The stages of the life cycle of primates vary considerably in duration. Among the most primitive members of the group, these stages are broadly comparable to those of other mammals of similar size. Higher in the phylogenetic scale they are substantially extended. The life span of a lemur is about 15 years, of a monkey 25–30 years, of a chimpanzee 40 years, and of man 70–75 years. The greatest differential between man and prosimians, indeed between the prosimians and apes and between apes and man, is in the duration of the infant and juvenile stages combined. The least differences occur in gestation period which, despite the general belief, can not be consistently correlated with adult body size. Gibbons, which weigh considerably less than macaques, have a 20 percent longer gestation period.

The clear trend toward prolongation of the period of juvenile and adolescent life is probably to be associated with the corresponding trend towards a progressive elaboration of the brain. The extended period of adolescence means that the young remain under adult (primarily maternal) surveillance for a long period, during which time the juvenile acquires, by example from its mother and peers, the knowledge that will allow it to become properly integrated as a fully adult member of a complicated social system. One might expect, therefore, a close correlation between the period of adolescence, the brain size, and the

Temperate
zone
species

complexity of the social system; and, insofar as the latter factor can be assessed, this appears to be the case.

Breeding periods. The reproductive events in the primate calendar are copulation, gestation, birth, and lactation. Owing to the long duration of the gestation period, these phases occupy the female primate (among higher primates anyway) for a full year or more; then the cycle starts again. The female does not usually come into estrus (physiologically receptive period) until the infant of the previous pregnancy is weaned.

The reproductive patterns of prosimians differ from those of anthropoids. Prosimians show one or more discrete breeding seasons during the year, during which time they may undergo more than one reproductive estrous cycle. The breeding seasons are separated by periods of anestrus, which in galagos are accompanied by changes in the skin of the vulva (the external organ of the genitalia), which closes over, completely sealing the vagina. In the Sudan estrus occurs only twice yearly in the African bush baby (*Galago senegalensis*), during December and August; but in captivity, breeding seasons may occur at any period in the year. In the wild, birth seasons are closely correlated with the prevailing climate, but in captivity under uniformly equable laboratory conditions, this consideration does not apply. In its native Madagascar the ring-tailed lemur (*Lemur catta*) has only a single breeding season during the year; conception occurring in autumn (April) and births taking place in late winter (August and September); but in zoos in the Northern Hemisphere a seasonal inversion of the birth period occurs in which it shifts to late spring and early summer. These examples suffice to indicate the influence of environmental factors on the timing of the birth seasons.

Reproductive cycles in the Anthropeida continue uninterrupted throughout the year, though seasonality in births is characteristic of primate species living either outside the equatorial belt (5° north and south of the Equator) or at high altitudes in equatorial regions, where dry seasons and seasonal food shortages occur. Seasonality of births in macaques (*Macaca* species) has been documented in Japan, on Cayo Santiago in the Caribbean (where an introduced population thrives under semi-natural conditions), and in India. Observations on langurs in India and Ceylon, on geladas in Ethiopia, and on patas monkeys in Uganda have also demonstrated seasonality in areas with well-marked wet and dry seasons. Anthropoids within the equatorial belt display birth peaks rather than birth seasons. A birth peak is a period of the year in which a high proportion of births, but not by any means all, are concentrated. Equatorial primates such as guenons, colobus monkeys, howlers, gibbons, chimpanzees, and gorillas might be expected to show a pattern of births uniformly distributed throughout the year, but population samples are as yet too small to make this assumption. In captivity, under both laboratory and zoo conditions, anthropoids breed throughout the year with little evidence of seasonality. Even in man there is evidence of high birth peaks. In Europe, the highest birth rates are reached in the first half of the year, and in countries in the Southern Hemisphere in the second half. This may, however, be a cultural rather than an ecological phenomenon, for marriages in certain western countries reach a peak in the closing weeks of the fiscal year, a fact which undoubtedly has some repercussions on the birth period.

Gestation period and parturition. The period during which the growing fetus is protected in the uterus is characterized by a considerable range of variation among primate species, but shows a general trend towards prolongation as one ascends the evolutionary scale. The range of prosimians overlaps that of the anthropoids, which is never shorter than 140 days. Tree shrews, for example, have a gestation period of 41–50 days, lemurs 132–134 days, macaque monkeys 146–186 days, gibbons 210 days, chimpanzees and gorillas 250–290 days, and man (on the average) 267 days. Even small primates like bushbabies have a gestation period considerably longer than those of nonprimate mammals of equivalent size, reflecting the increased complexity and differentiation of primate structure compared with that of nonprimates. Although in primates

there is a general trend towards evolutionary increase in body size, there is no absolute correlation between body size and the duration of the gestation period. Marmosets, for example, are considerably smaller than spider monkeys and howler monkeys but have a slightly longer pregnancy (howler monkeys 139 days; marmosets 140–150 days).

An extraordinary and somewhat inexplicable difference exists between the dimensions of the pelvic cavity and the dimensions of the head of the infant at birth in monkeys and man on the one hand, apes on the other. The head of the infant ape is considerably smaller than the pelvic cavity, so that birth occurs easily and without prolonged labor. When the head of the infant monkey engages in the pelvis the fit is exact, and labour may be a prolonged and difficult affair as it is in man, in whom similar conditions prevail. Human parturition, however, is generally a much more extended process than that of monkeys. Like the human infant, the monkey is born head first. Twin births are as rare in Old World monkeys and apes as they are in man, but certain New World monkeys (e.g., marmosets) habitually produce twins.

The degrees of maturation and mother dependency at birth are obviously closely related phenomena. Newborn primate infants are neither as helpless as kittens, puppies, or rats, nor as precocial as the newborn gazelles, horses, and other savanna-living animals. With a few exceptions, primate young are born with their eyes open and are fully furred. Exceptions are the lowly tree shrews and mouse lemurs, which carry their young in their mouths in the manner of many mammals which bear altricial infants. Primate life being peripatetic, it is axiomatic that the infants must be able to cling to the mother's fur. The young of most higher primates have fully prehensile hands and feet at birth and are thus able to grasp the maternal fur without assistance; only man, chimpanzee, and gorilla need to support their newborn infants. In man this stage persists throughout the nursing period (six months or more), but in gorillas it lasts only four weeks and in chimpanzees for an even shorter period.

It seems likely that the difference between the African apes and man in respect of postnatal grasping ability is related to the acquisition in man of bipedal walking. One of the anatomical correlates of the human gait is the loss of the grasping function of big toe, which is aligned in parallel with the remaining digits. Such an arrangement precludes the use of the foot as a grasping extremity. The human infant and, to a lesser degree, the gorilla infant must depend largely on its grasping hands to support itself unaided. The fact that man is habitually bipedal and that, consequently his hands are freed from locomotor chores, may also be a contributory factor: the human mother can move about and at the same time continue to support her infant. Selection for postnatal grasping, therefore, has not had the high survival value that it has in nonhuman primates in which the survival of the infant depends on its ability to hold on tightly. It is a well-known fact, of course, that newborn human infants can support their own weight, for short periods, by means of their grasping hands. Clearly adaptations for survival are not wholly lacking in the human species. Perhaps there are involved cultural factors which have the effect of suppressing natural selection for early infant grasping ability. The first may be that the social evolution of a division of labour between the sexes and a fixed home base has allowed the mother to park her infant with other members of the family as baby-sitters; the second that in more peripatetic communities, the invention of infant-carrying devices, such as the papoose technique of North American Indians, has made it unnecessary for the infant to support itself. Whatever the biological or cultural reasons, the human infant is less precocial than the young of any other primates.

Once the primate infant learns to support itself by standing on its own two (or four) feet, the physical phase of dependency is over; the next phase of psychological dependency lasts much longer. Man, metaphorically, is tied to his mother's apron strings for much longer periods than are the nonhuman primates. The reasons for this are discussed below. According to Adolph Schultz, the anthropologist whose comparative anatomical studies have

Condition of the newborn primate

Distribution of births through the year

illuminated knowledge of nonhuman primates during the last 40 years, the juvenile period of psychological maternal dependency is 2½ years in lemurs, 6 years in monkeys, 7–8 years in apes, and 14 years in man.

Growth and longevity. The characteristic pattern of weight-growth curves of man is also shared by higher nonhuman primates such as chimpanzees (*Pan troglodytes*) and rhesus monkeys (*Macaca mulatta*). There seems little reason to doubt that similar curves will be found to exist in other higher primates not yet studied. It has been suggested that the growth curves of primates are closely akin to the basic curves of mammals in general.

The trend toward the prolongation of postnatal life as a characteristic of primates affects all postnatal life periods including infantile, juvenile, adult, and the period of senescence. The period which shows the least change in apes and man is that of gestation. The relative stability of the gestation period is a clear indication that, developmentally speaking, in the degree of structural complexity and cellular differentiation, there is little difference between apes and man. The differences that undoubtedly exist can probably be attributed to cultural factors. There is no reason to suppose that apes are not built to last equally as long as man; it is simply that the natural hazards of gorilla or chimpanzee life are not buffered by the blessings of civilization, except under conditions of captivity, in which case life spans of 40 years and more are known. The potential life span of the chimpanzee has been estimated at 60 years.

The characteristic growth spurts of human infants in weight and height also occur in nonhuman primates but start earlier in the postnatal period and are of shorter duration. Primates differ from most nonprimate mammals by virtue of a delayed puberty in both sexes until growth is nearly complete; in male macaque monkeys, as in man, the growth spurt precedes the onset of puberty. To postpone this latter event is of selective advantage because it ensures that the juvenile males in a primate society are not sexually competitive with adult males until they are large enough and strong enough to contribute fully to the responsibilities of the society in which they live by taking part, for instance, in defense against predators.

Locomotion. Primate locomotion, being an aspect of behaviour which arises out of anatomical structure, shows much of the conservatism and the opportunism that characterizes the primates. Primates with remarkably few changes in their skeletons and musculature have adopted a bewildering variety of locomotor patterns. The "natural" habitat of primates—in the historical sense—is the canopy of the forest. Although many primates have adopted the ground as their principal foraging area during the day, given the opportunity, they return to the trees to sleep at night. Trees provide cover from the climate and protection from predators; they are of course also a source of food. Only the gelada (*Theropithecus gelada*), the hamadryas

baboon (*Papio hamadryas*) of the mountainous regions of Ethiopia, and the chacma baboon (*P. ursinus*), which lives on the rocky coast of the Cape of Good Hope, South Africa, are ground sleepers; even these animals seek the protection of the cliffs and rocky precipices of their habitats at night. No primate sleeps totally unprotected; as a consequence of their relative immunity from predation, primates are heavy sleepers.

The essential arboreality of primates has guaranteed the relative uniformity of the locomotor apparatus. Even man, who has long since abandoned the trees as his principal lodging place, has only partially lost the physical adaptations for tree climbing; his hands remain in the arboreal mold. Only his feet have lost their primitive prehensility in adapting to bipedal walking. Primate locomotion can be classified on behavioral grounds into four major types (see Table 1). Within these major categories there are a number of subtypes, and within these subtypes there are an infinite number of variations between species and, by virtue of individual variability, within species. The differences between the four major categories lie principally in the degree to which the forelimbs and hindlimbs are used to climb, swing, jump, and run.

Vertical clinging and leaping, for instance, is primarily a function of the hindlimbs, as is bipedalism, whereas brachiation is performed exclusively with the forelimbs. Quadrupedalism involves both forelimbs and hindlimbs, of course, although not to an equal extent. Some quadrupeds are hindlimb dominated; in others, the forelimb and the hindlimb are equally important. The hindlimb-dominated primates, such as the langurs and colobus monkeys, employ a large element of leaping in their movements, a less notable feature of the more generalized quadrupeds such as the species of the genus *Cercopithecus*. The quadrupedal category is inevitably somewhat of a grab bag, and the gaits included in it have not yet been studied critically. One subtype, here designated as slow climbing, differs profoundly from the other subtypes of the category, being somewhat ponderous and devoid of elements of leaping or jumping. The species in this category are all arboreal nocturnal prosimians.

As many authorities who have studied locomotion in free-ranging primate species have pointed out, the classifications of locomotion into categories is a somewhat artificial procedure. A chimpanzee shows a variety of different gaits according to the circumstances of the environment, and brachiation may form only a minor part of its repertoire; this holds true also for the langurs and colobus monkeys, which are placed in the subtype Old World semibrachiation. Although the categories are phrased in behavioral terms, their implications are strictly anatomical. Brachiation is the mode of locomotion for which the animal is specifically adapted; the anatomical correlates of brachiation are quite unmistakable and can be determined equally as well in fossil bones as in living animals. A

Locomotor types

Average and potential life span

Table 1: Locomotor Classification

category	subtype	activity	primate genera
Vertical clinging and leaping		leaping in trees and hopping on the ground	<i>Avahi, Galago, Hapalemur, Lepilemur, Propithecus, Indri, Tarsius</i>
Quadrupedalism	slow climbing type	cautious climbing—no leaping or branch running	<i>Arctocebus, Loris, Nycticebus, Perodicticus</i>
	branch running and walking type	climbing, springing, branch running and jumping	<i>Aotus, Cacaiao, Callicebus, Callimico, Callithrix, Cebuella, Cebus, Cercopithecus, Cheirogaleus, Chiropotes, Lemur, Leontideus, Phaner, Pithecia, Saguinus, Saimiri, Tupaia</i>
	ground running and walking type	climbing, ground running	<i>Macaca, Mandrillus, Papio, Theropithecus, Erythrocebus</i>
	New World semi-brachiation type	arm swinging with use of prehensile tail; little leaping	<i>Alouatta, Ateles, Brachyteles, Lagothrix</i>
Brachiation	Old World semi-brachiation type	arm swinging and leaping	<i>Colobus, Nasalis, Presbytis, Pygathrix, Rhinopithecus, Simias</i>
	true brachiation modified brachiation	gibbon type of brachiation chimpanzee and orangutan type of brachiation	<i>Hylobates, Symphalangus Gorilla, Pan, Pongo</i>
Bipedalism		striding	<i>Homo</i>

Source: From J.R. and P.H. Napier, *A Handbook of Living Primates*.

purely behavioral classification, based on observations in the living, can have little value in palaeontological diagnosis. It is quite arguable that although pure arm swinging is rare in Old World semibrachiators, it is nevertheless one of the principal survival mechanisms of the group, an adaptive advantage that in critical circumstances provides these species with an "edge" over their competitors. In some instances it may well be the case that the particular category is of evolutionary significance only. The gorillas, for instance, whose anatomy is that of a brachiator, in fact, seldom brachiate; adult males probably never do. Gorillas are large animals and bulk-food eaters, forced by their dietary demands to obtain their food on the ground; living gorillas have adapted to a terrestrial life by adopting a compromise gait in which they walk quadrupedally bearing their weight on the backs of their knuckles. This is a form of quadrupedalism (inasmuch as all four limbs are employed) but, once again, the anatomy of their bodies indicates that in past time gorillas (or their evolutionary forebears) were smaller animals which lived and moved in trees. The modern gorilla differs widely from its historical counterpart in behaviour but, structurally, its kinship with its arboreal ancestor is overtly apparent.

Changes in climate and geography during the evolutionary history of primates may also have led to structural atavisms in the anatomy of living primates. Many chimpanzees now living in woodland-savanna conditions in Africa, where the trees are widely spaced and generally unsuitable for the classic arm swinging progress of a brachiator, have adopted a largely ground-living life. Like the gorilla, the chimpanzees are also knuckle walkers, but given an environment like that of a zoo with a cage specially designed with a plentitude of overhead bars and ropes, they brachiate freely and frequently. Habitat changes in the past undoubtedly account for many apparent anomalies in the locomotor classification of present-day primates.

When the subject of primate arboreal locomotion is studied in evolutionary terms, through the medium of fossils, it becomes clear that locomotor categories are not discrete, but that they constitute a continuum of change from a hindlimb-dominated gait to a forelimb-dominated one. The best single indicator of gait, one that has the added advantage of being strictly quantitative, is the intermembral index. Briefly, the index is a ratio expressed as percentage of arm length to leg length; an index over 100 indicates relatively long arms. The intermembral index provides a model by means of which the locomotion of an early primate can be inferred by determination of the intermembral index of the fossil skeleton. Animals do not necessarily fall discretely into categories. Species with indices lying between those of clearly recognizable locomotor types represent transitional types, whose style of locomotion manifests features of both of the bracketing categories. Some lemurs have indices that fall between 65 and 75, but their gait is a combination of vertical clinging and quadrupedalism. The ring-tailed lemur in fact shows an intermediate form of behaviour; it is generally quadrupedal but, given the right environmental milieu, can be observed in vertical clinging and leaping. The South American spider monkeys (*Ateles*), whose index lies between 100 and 108, show a type of locomotion that contains the elements of both quadrupedalism and brachiation.

Applying the intermembral index to fossil primates, it has become apparent that the earliest primates living in the Eocene Epoch, from about 54,000,000 to 38,000,000 years ago, moved about in the manner of modern vertical clingers and leapers. Quadrupedalism was well-established during the Miocene Epoch (about 26,000,000 to 17,000,000 years ago) when the two major environmental types of quadrupedal gait—the terrestrial and the arboreal—were established, with indices in the region of 85–100 and 75–85, respectively. Brachiation, associated with a high intermembral index, was established as a way of arboreal life at the end of the Miocene and the beginning of the Pliocene.

The origin of bipedalism, which implies the appearance of manlike creatures, is indeterminate. There is direct evidence of bipedalism extending back 1,500,000

to 2,000,000 years, and certain indirect evidence (see below *Evolution and paleontology*) suggests that bipedalism might have evolved in a modified form some 10,000,000 to 14,000,000 years ago.

Locomotion of primates is an area of primate biology of great interest to scholars. It is not only that locomotion is a basic factor of primate life (primates must move to eat and must eat to live) but it is also an indicator of evolutionary progress. To study the origins of primate locomotion is to study primate evolution. Because of its immediacy for survival, natural selection has acted strongly on locomotor behaviour and the locomotor apparatus; changes in locomotor pattern may thus be said to provide principal milestones on the road to modern primates.

Bipedalism, far from being a unique possession of mankind, as is generally thought, is a basic possession of the order Primates. With one or two exceptions (one of which is the tree shrew, which may not even be a primate), all primates sit upright. Many stand upright without supporting their body weight by their arms, and some actually walk upright for short periods. The view that the possession of uprightness is a solely human attribute is untenable; man is merely one species among the 189 that constitute the order who has exploited the potential of his ancestry to his ineffable advantage.

Chimpanzees, gorillas and gibbons, macaques, spider monkeys, capuchins, and others are all frequent bipedal walkers. To define man taxonomically as "bipedal" is not enough; to describe him as habitually bipedal is nearer the truth, but habit as such does not leave its mark on fossil bones. Some more precise definition is needed. Man's walk has been described as striding, a mode of locomotion defining a special pattern of behaviour and a special morphology. Striding, in a sense, is the quintessence of bipedalism; it is a means of travelling during which the energy output of the body is reduced to a physiological minimum by the smooth, undulating flow of the progression. It is a complex activity involving the joints and muscles of the whole body in its performance and it is likely that the evolution in the human gait took place gradually over a period of 10,000,000 years or so.

The pattern of locomotion of man's ancestors immediately preceding the acquisition of bipedalism has long been a matter of controversy, and the question has not yet been resolved. The evidence derived from anatomical, physiological, and biochemical studies for the close affinity of chimpanzees and gorillas with man would imply that he evolved from a brachiating ancestry. But the many differences between the apes and man, particularly in the structure of the limbs and teeth, have persuaded the majority of leading anthropologists that the theory of a brachiating origin of man is untenable. The issue is still hotly debated, the more so since current immunological studies of serum proteins are tending to emphasize the recentness of the separation of pongid and hominid stocks (the lines that led, respectively, to apes and man) in phylogeny. An alternative to the "brachiating theory" must inevitably be some sort of "quadrupedal theory." Various authorities have proposed other solutions: semibrachiation (a form of quadrupedalism with emphasis on arm swinging), a formal quadrupedal gait, and even a form of locomotion similar to that of *Tarsius* and other clingers and leapers. At the present time, there is insufficient information to elucidate the phylogeny of man's bipedal gait. It can, however, be assumed that it must have involved a large measure of truncal uprightiness, a good deal of leaping, and a sizable element of arm swinging. There is no living primate that demonstrates, precisely, this combination of locomotor characteristics, although the arboreal monkeys of the subfamily Colobinae theoretically fit the bill better than most.

Ecology and diet. Few primate species are found today in temperate latitudes. The limiting factors appear to be the general climate and the duration of the daylight hours. The higher the latitude, the greater is the fluctuation between summer and winter temperatures and the shorter are the hours of winter daylight. The combination of quiescent growth periods of fruits, leaves, and grasses, a frozen top soil, and a restricted diurnal period all make

The relation of locomotion to habitat changes

The intermembral index

Walking

conditions extremely hard for monkeys, which depend on foraging in the trees, on the ground, or under it for their food supply. No monkey lives, or has ever lived, in the zones of permafrost, but it is probable that in times past they flourished in temperate regions. Without the harassment of man, monkeys no doubt could have survived in the temperate zones; their foraging, albeit condensed into a temperate-zone day of seven hours or so, would have supplied them with the food they required, providing they were free to forage without the threat of human persecution. Dietary adaptability of primates is greater than realized. On the Shimokita Peninsula of the Japanese Island of Honshu and in the Shiga Mountains of the same island, Japanese macaques survive from late November until April on incredibly meagre rations. The diet of these monkeys, during the months when snow covers the ground and when the last remaining summer seeds and leaves have vanished, is bark and the fleshy cambium layer that lies beneath it. The Japanese macaques, cold and hungry, ride out the worst that winter, at the latitude of 41° N, offers. Near Moscow, several species of African primates have been kept for several years in outdoor facilities in which winter temperatures go well below -18° C (0° F). It is true that they are well provisioned but, nevertheless, it is almost incredible that a tropical species should be sufficiently adaptable and possess adequately effective homeostatic mechanisms to survive in this degree of climatic adversity. Ruggedness of this sort is a primate characteristic, a survival mechanism that operates in the worst conditions, and one which provides a clue to the real reason for the primacy of the primates among mammals.

Diet. The diet of primates is a factor of their ecology that, during their evolution, has clearly played an important role in their dispersion and adaptive radiation as well as in the development of the teeth, jaws, and digestive system. Diet is also closely related to locomotor pattern and to body size.

The principal food substances taken by primates may be divided into vegetable (fruits, flowers, leaves, nuts, barks, pith, seeds, grasses, stems, roots, and tubers) and animal (birds' eggs, lizards, small rodents and bats, insects, frogs, and crustacea). The flesh of larger mammals (including primates) is not listed as an important item of nonhuman primate diet, although it is taken by certain species such as baboons and chimpanzees in special circumstances that are not yet fully understood.

While in many mammalian groups, diet is selective and specific to the order, among primates it is difficult to establish any hard and fast rules. Although there are decided preferences for certain food items, catholicity is more characteristic than specificity. Generally speaking, primates are omnivorous, as the physiology of their digestive system attests. Relatively few examples of dietary specialization are to be found. The so-called leaf-eating monkeys, a sobriquet that embraces the whole of the subfamily Colobinae, are by no means exclusively leaf eaters and according to season include flowers, fruit, and seeds in their diet. The howler monkeys of the New World have a similar dietary preference.

Broadly, however, certain overall dietary preferences are discernible. The leaf-eating langurs have already been mentioned. The apes (other than the gorilla) are substantially fruit eaters; the percentages of fruit in the diet of the gibbon and the three great apes are as follows: gibbon 62 percent, orang 50 percent, chimpanzee 67 percent, and gorilla 2-15 percent.

Many of the smaller nocturnal prosimians such as tarsiers, galagos, dwarf lemurs, sportive lemurs, the aye-aye, slender lorises, and the South American night monkey are substantially insectivorous.

The larger diurnal prosimians (e.g., typical lemurs, the sifaka, and the indri) and the large nocturnal prosimians (the Asian slow loris and the African potto) are more vegetarian, including both fruit and leaves in their menu. It seems apparent that size, rather than activity rhythm, governs the nature of the primate diet. The small marmosets of the South American genus *Callithrix* have exclusively diurnal rhythms and are insectivorous, while the slightly

larger, but equally diurnal, tamarins (*Saguinus*) are more omnivorous.

The most widespread monkey populations are those of the subfamily Cercopithecinae. Many ground-living species are included in this group, the members of which are notable for the extreme catholicity of their diet; meat eating as a regular pattern, however, has not been reported among them.

Size in relation to food habits. In evolutionary terms, increase in size has probably played a large part in determining the direction of primate evolution. Fossil evidence of early primates of the Paleocene (about 65,000,000 years ago) suggests that these were small, forest-living creatures, not yet specialized as tree climbers. Their bodily form suggests that they possessed the build and the gait of squirrels. Their molar teeth bore high, pointed cusps but were neither as tall nor as pointed as those of the insectivore-like ancestors of the preceding Cretaceous epoch, whose molars were ideally adapted for cracking the hard, chitinous exoskeletons of insects. This fact suggests that the reduction of the molar cusps was associated with the adoption of a fruit-eating habit. Although this has some validity as a generalization, it should not be taken too literally. All primates include insects in their diet and there are many almost exclusively insectivorous forms, which have reduced the height and acuity of their molar cusps. Increasing body size, a trend which is clearly apparent throughout primate evolution, would have been associated with the adoption of supplementary sources of food. Insect eating is a laborious way for a large animal to obtain nourishment unless it happens to possess the extensible tongue of those mammalian "vacuum cleaners," the anteaters.

An increase in size and the gradual addition of bulk foods to the diet would, in turn, have affected the habitat and the pattern of locomotion of primates. Fruit eating is more easily accomplished in trees than on the ground and this growing preference amongst primates may have been accompanied by the tree-climbing habit. Once such a habit was established, natural selection would have operated to espouse the suitable physical adaptations which would facilitate such a behaviour. Suitable adaptations in this case would have been the facility to climb, leap, and balance in trees.

It is noteworthy that during evolution the development of a prehensile foot preceded that of a prehensile hand. Vertical clinging primates such as the tarsiers or small squirrel-like quadrupeds like the marmosets and tamarins—all of which have prehensile feet, but not prehensile hands—by remaining small, have never been subject to the same evo-

From J.R. Napier and P.H. Napier, *A Handbook of Living Primates* (1967), Academic Press

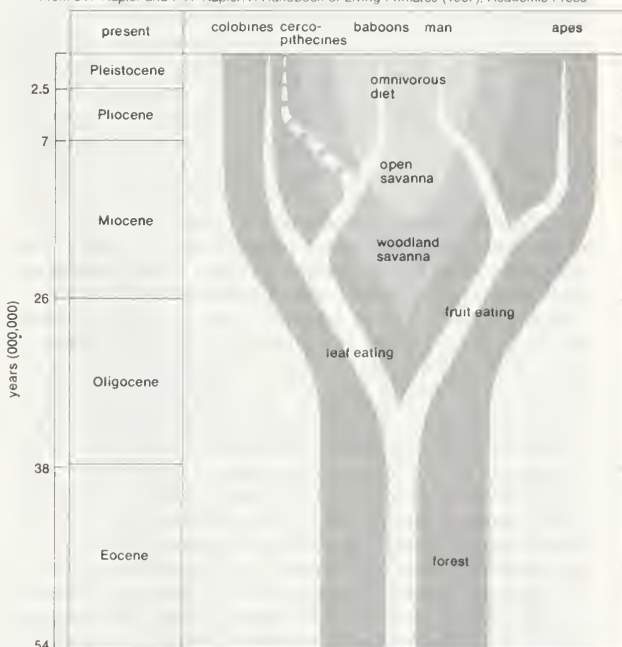


Figure 21: Hypothetical evolution of primate feeding habits in relation to vegetational areas.

Resistance to cold

The shift from insect to fruit eating

lutionary pressures that have impinged on larger primates. A large arboreal primate without prehensile hands is at a considerable disadvantage in moving about in the canopy of trees, but a small one suffers little disadvantage. Amid the large and firm branches, size is no particular hazard, but at the periphery of the crown, where the fruit is most abundant and the branches are slender and flexible, the risk of falling is increased. It is, therefore, likely that the combination of an increase in body size, associated with the inevitable shift towards a bulk diet, led first to the evolution of a grasping hand, then to the appearance of a prehensile hand, and finally to an opposable thumb. Four prehensile extremities are more effective than two in defying gravity.

Such adaptations of the forelimbs would have had the effect of equalizing the role of the limbs. The limbs of vertical clingers are functionally disparate, the lower pair being dominantly propulsive, the upper, secondary and purely supportive. The limbs of quadrupeds, on the other hand, are more homogeneous, both pairs having a propulsive function during running. Thus, it would seem that the transition in locomotor grade between vertical clinging and leaping and quadrupedalism came about as adaptation to increased body size. Size, diet, ecology, locomotion, and anatomical structure provide a constellation of causes and effects that are critical factors in the evolution of the primates.

Effect of forest stratification. An aspect of primate ecology that has not yet been fully investigated is the effect of the stratification of forests on the adaptations (particularly the locomotor adaptations) and on the speciation of primates. In recent years it has become apparent that forests are not simply a random consociation of trees but are an ordered system of plants under the control of climatic, photobiotic, and edaphic (soil) factors.

The chief physiognomic features of rain forests, the ancestral home of primates and the principal habitat (in numerical terms) of nonhuman primates today, are the evergreen broad-leaved trees that, collectively, form a closed canopy, so opaque to sunlight that the forest floor is in perpetual twilight. Epiphytes and thick-stemmed lianas drape the trees, linking one crown to another and providing aerial pathways for monkeys to pass from tree to tree through a continuum of interlacing branches, a three-dimensional maze that provides home, restaurant, shopping districts, and highways for primates. Three strata of the canopy of rain forests are usually recognized by botanists: an understory, a middle story, and an upper story. The understory is often "closed," the crowns of the constituent trees overlapping one another to form a dense, continuous, horizontal layer. The middle story is characterized by trees that are in lateral contact but do not overlap; and the highest story by tall trees, some 50 metres (160 feet) or more, which form a discontinuous layer of umbrella-shaped crowns. The occasional "emergent" forest giant may tower above the highest layer of the canopy. There is some evidence, much of it conflicting, that zonation of forest primates occurs within the forest canopy. The stratification of forest is extremely variable and the number of layers tends to diminish from three to two in secondary forest, dry deciduous forest, and montane forest; and from two to one as temperate zone, tropical woodland, or montane woodland supervenes.

Tropical grasslands, or savannas, are also the homes of primates in Africa and Asia; no savanna-living primates exist in South America. Tropical grasslands comprise a mixture of trees and grasses. The proportion of trees to grass varies directly with the rainfall; areas of high seasonal rainfall support single story woodlands of tall trees while lush grasses form the ground vegetation; but where rainfall is both seasonal and low, the trees consist of stubby xerophilous (dry-loving) shrubs and short, tussocky grasses. The principal primate fauna of the savanna biome are the ground-living species; in Africa, the grivet (*Cercopithecus aethiops*), baboons (*Papio* species), and the patas monkey (*Erythrocebus patas*); and in Asia the macaques (*Macaca* species) and the Hanuman langur (*Presbytis entellus*).

Tropical montane forests or tropical rain forests at high altitude also abound in primates in Africa, Asia, and South

America. In equatorial Africa certain primate species have colonized the montane-savanna regions, or moorlands, where the rugged mountainous terrain and seasonal food scarcity support herds of geladas and hamadryas baboons. These high mountaineers of Africa have no ecological counterparts in Asia or South America.

As a measure of the physiological adaptability of primates it may be noted that various populations have at different times been subjected to relocation in alien environments by man. Some of these expatriates are the green monkeys (*Cercopithecus sabaeus*) of St. Kitts in the West Indies, the mona monkeys (*C. mona*) of Grenada, the macaques on Mauritius, and the Barbary apes (*Macaca sylvana*) of Gibraltar.

FORM AND FUNCTION

Distinguishing features. The basis of the success of the order Primates is the relatively unspecialized nature of their structure and the highly specialized plasticity of their behaviour. This combination has permitted the primates throughout their evolutionary history to exploit the wide variety of novel ecological opportunities that have come their way. Although there are a few highly specialized species among the lower primates (the aye-aye, the tarsier, the potto, and the lorises, for instance), the higher primates, the anthropoids, are extremely conservative in their structure; morphologically speaking, they have maintained a position in the evolutionary midstream and have avoided the potential stagnation of specialized life near the bank. Specialization is not always a liability; in times of environmental stability the specialized animal enjoys many advantages, but in a rapidly changing world it is the less specialized animals that are more likely to survive and flourish.

The plasticity of primate behaviour is largely a function of the brain. The primate brain is distinguished by its relatively large size compared with the size of the body as a whole; it is also notable for the complexity and elaboration of the cerebral cortex, the function of which is to receive, analyze, and synthesize the incoming impulses from the sense organs and to convert them into appropriate motor actions, which in turn constitute behaviour.

It has been observed that there is no distinguishing feature which characterizes all primates except a negative one—their lack of specialization. St. George Mivart defined the order Primates in the following terms:

... an unguiculate, clavicate, placental mammal with orbits encircled by bone; three kinds of teeth at least at one time of life; brain always with a posterior lobe and a calcarine fissure; the innermost digits of at least one pair of extremities opposable; hallux with a flat nail or none; a well-developed caecum; penis pendulous; testes scrotal; always two pectoral mammae.

When these criteria are examined one by one the inevitable conclusion is reached that, singly, none of these characters is unique to primates nor do all primates demonstrate every single character. Collectively, however, more primates than nonprimates show this particular constellation of characters. Perhaps a more meaningful way to look at primate characteristics is in the form of evolutionary trends. This approach is quantitative, inasmuch as it is recognized that not all adaptations are developed in primates to the same degree. A list of evolutionary trends with behavioral connotations is also a somewhat more functional inventory than Mivart's essentially morphological one and, furthermore, is more applicable to the study of fossil primates.

Primates are essentially arboreal animals whose limbs are adapted for climbing, leaping, and running in trees. Active arboreal life requires the mechanical assistance of a long tail and sensitive, grasping hands and feet with opposable thumbs and big toes, to aid in climbing and to ensure stability on slender branches high above the ground. Active arboreal locomotion also requires a much more accurate judgment of distances than life on the ground; this is facilitated by the development of stereoscopic vision, the anatomical basis of visual judgments in depth. The forward-facing eyes of primates are adaptations for this type of visual precision. A highly developed sense of smell is not nearly as important for animals leading

Forest structure and traffic ways

The importance of a large brain

Binocular vision

an arboreal life as it is for those on the ground. Primates thus have a much reduced olfactory mechanism; noses are shorter, the turbinal (or scroll) bones of the nose are reduced in number and complexity compared with most nonprimate mammals.

Primates, being on the whole omnivorous animals, have evolved a multipurpose dentition; incisors to cut, canines to pierce and tear, and molars and premolars to crush and grind. Such a dental arrangement equips its possessor with the power to cope with many different sorts of food—insect food, vegetable food, or flesh—and provides a major contribution for the generalized way of life that is so characteristic of the order.

A basic adaptation for the primate way of life (including the unique specialization of man—upright walking) is verticality of the trunk. Above the level of the tree shrews, the ability to hold the body upright is a universal primate possession; indeed paleontological evidence indicates that truncl uprightness was a primate characteristic as far back as the early Eocene, some 54,000,000 years ago.

Above all, the principal evolutionary trend of primates has been the elaboration of the brain, particularly of that portion of the cerebral hemispheres known as the neopallium or neocortex. A neocortex is characteristic of higher vertebrates, such as mammals, which operate under the control of multiple sources of sensory input. In many mammals the olfactory system dominates the senses and the cerebral hemispheres consist largely of paleocortex—the “smell brain”—of lower vertebrates. The arboreal habit of primates has led to a dethronement of the olfactory sense and the accession of a tactile, visually dominant sensory system. This evolutionary trend has resulted in the dramatic expansion and differentiation of the neocortex.

General structure. *Vertebral column and posture.* All primates retain a clavicle or collarbone (lost in many mammalian groups), a separate radius and ulna in the forearm, and a separate tibia and fibula in the lower leg. The single exception to this is the tarsier, in which the fibula becomes fused to the tibia.

The primate vertebral column shows a basic mammalian pattern of components including an “anticlinal” vertebra, situated in the mid-thoracic (upper-back) region of the spinal column and marking the transition between the forelimb and hindlimb segments. In a galloping greyhound the anticlinal vertebra is at the apex of the acute curve of the back. An anticlinal vertebra is characteristic of all quadrupeds and is seen in all primate families except the hominids and the pongids, whose posture is upright or semi-upright. The evolutionary trend in the vertebral column is toward shortening the lumbar, sacral, and caudal regions. Extreme shortening in the lumbar region, with complete loss of the caudal (tail) vertebrae, occurs in gibbons, the great apes, and man. Prehensibility of the tail is a specialization of certain New World monkeys, but appears also for a brief period in the infants of many Old World monkey groups, in which it provides an important mechanical aid to survival.

Hands and feet. With two exceptions, all primates have retained five digits on hand and foot. The exceptions are the spider monkeys of South America and the colobus monkeys of Africa, which have lost or reduced the thumb. This appears to be an adaptation for locomotion, the rationale for which is not fully understood at present.

All (except tree shrews) possess prehensile (grasping) hands, and (except man) prehensile feet. The hands of Anthropoidea show a greater range of precise manipulative activity than those of Prosimii. Prosimians lack the functional duality of the hands of anthropoids. Duality in hand function has been described for man and nonhuman primates in terms of precision and power grips. The power grip of prosimians is very well developed but the precision grip is lacking. Among Anthropoidea, the New World monkeys show a considerable advance over prosimians in tactile sensitivity but they possess less functionally effective hands in prehensile terms than Old World monkeys. The critical component of the prehensile hand in terms of skilled manipulation is the opposable thumb; a thumb, that is to say, that is capable of being moved freely and independently. The movement of opposition is a rotary

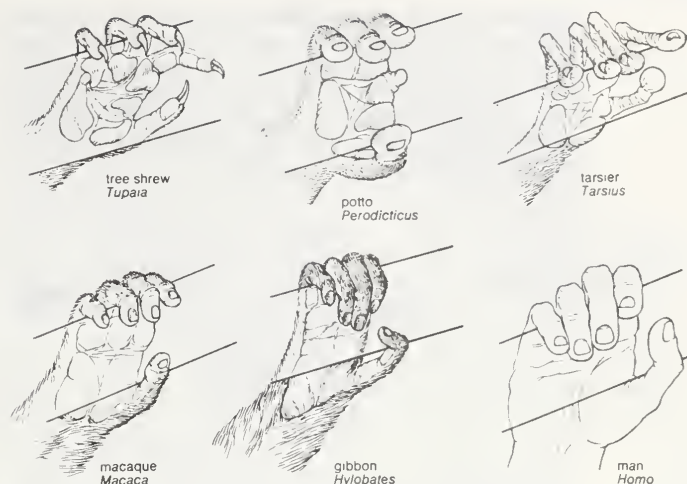


Figure 22: Shape and function of the hand in certain primates.

Drawing by R. Keane

movement in which the thumb, swinging about its own axis, comes to face the under or palmar surface of the tips of the fingers. The opposable thumb is the basis of the precision grip which, though present to some extent in all primates, is particularly highly developed in man. Opposability is present to some degree or other in most primates but varies considerably in its functional effectiveness as an instrument of fine manipulation. The baboons and man are pre-eminent in this respect. The apes, having short thumbs and long fingers, are handicapped in relation to delicate manual dexterity but are adept in the coarser elements of hand use, particularly in relation to tree climbing and branch swinging.

Teeth. A dentition with different kinds of teeth (heterodonty)—incisors, canines, and cheek teeth—is characteristic of all primates and indeed of mammals generally. Heterodonty is a primitive characteristic and primates have evolved less far from the original pattern than most mammals. The principal changes are a reduction in the number of teeth and an elaboration of the cusp pattern of the molars.

The primitive mammalian dental formula is assumed to have been $\frac{3 \cdot 1 \cdot 4 \cdot 3}{3 \cdot 1 \cdot 4 \cdot 3} = 44$ teeth (the numbers being the

numbers respectively of pairs of incisors, canines, premolars, and molars in the upper and lower jaws). No living primate has retained more than two incisors in the upper jaw and only in tree shrews is the primitive number of three incisors retained in the lower jaw. The incisors are subject to considerable variation in prosimians. Characteristically, the upper incisors are peglike, one or other pair often being absent; in the lower jaw the incisors show a peculiar conformation which has been likened structurally and functionally to a comb. The dental “comb” of prosimians is composed of the lower canines and lower incisors compressed from side to side and slanted forwards; the dental comb appears to be primarily a device for grooming the fur. Canines are present throughout the order, but show remarkable variation in size, shape, projection, and function. Characteristically the teeth of cercopithecines have a function in the maintenance of social order within the group as well as an overtly offensive role; their function as organs of digestion is relatively unimportant. They are large and subject to sexual dimorphism, being larger in males than females. Man’s canines are small and are almost wholly dental in function, with no size differences between the sexes.

The trend in the evolution of the cheek teeth has been to increase the number of cusps and reduce the number of teeth. Both molars and premolars show this tendency. No living primate has four premolars; prosimians and New World monkeys have retained three on each side of each jaw, but at the anthropoid level two premolars are the invariable rule. The primitive premolars are uniform in shape and are unicuspid, but in primates the most posterior premolar tends to evolve either one or two extra

Evolutionary reduction of the cheek teeth

Power versus precision grips

cusps (molarization), an adaptation providing an extended tooth row for a herbivorous diet. In species with large upper canines, the most anterior lower premolar assumes a peculiar shape known as sectorial, functioning as a hone for the scythelike canine. In man, whose canines are small and unremarkable, the first and second premolars are identical in shape and two-cusped.

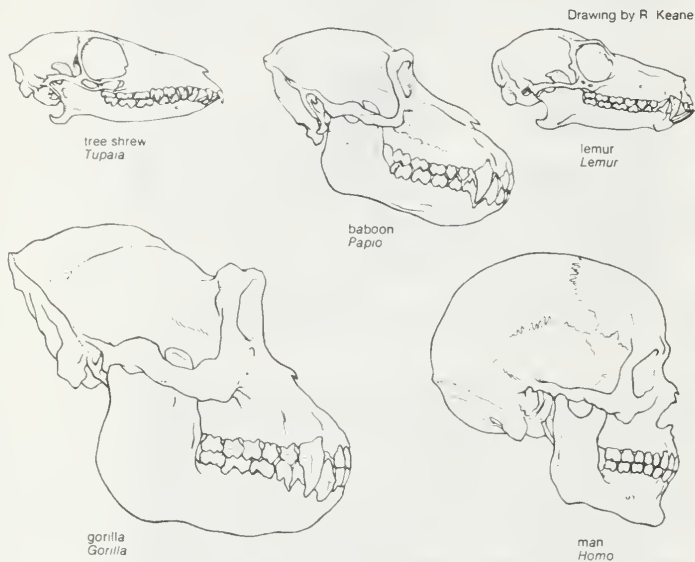


Figure 23: Skulls and teeth of primates.

The trend in the morphology of the molars has been to increase the primitive three cusps to four or five. The apes and man are more specialized in this feature than are the monkeys or the lower primates, having four cusps on the molar crown in the upper jaw and five cusps on the lower. A further tendency has been to reduce the molar series from three to two in both jaws, principally in the lower.

Snouts, muzzles, and noses. The reduction of the snout in primates is a correlate of the diminution of the olfactory sense (see below *Sensory reception*). To a great extent visual acuity and manual dexterity have replaced the sensitive, inquiring nose found in so many nonprimate mammals. A marked reduction in the complexity of the "scroll" bones of the nose, the richness of the innervation of the olfactory mucous membrane, and the sensitivity of the moist tip of the nose—the rhinarium—is associated with the reduction in length of the primate snout. Although the trend in primate evolution is towards a dethronement of the primacy of the sense of smell, there are still some good snouts to be seen in those lower primates, that retain a naked moist rhinarium attached to the upper lip.

Prosimians depend for many aspects of their social and reproductive behaviour on olfactory signals which are effected by means of special scent glands distributed in different regions of the body but congregated principally in the anal and perineal regions. Marking behaviour, the placing of scent at various points in the environment, is a prominent feature of the prosimian repertoire of communication. Marking behaviour ceases to be of much importance above the prosimian level, and in anthropoids the structures by which olfactory signals are given and received are diminished. Monkeys, apes, and man, though they may be lower down the scale of olfactory excellence than prosimians, rely in considerable measure on the sense of smell; all higher primates, for instance, including man, sniff at unfamiliar items of food before placing them in the mouth.

The shape of the nose of higher primates is one of the most reliable means of distinguishing Old World monkeys from New World monkeys at a glance. In New World monkeys or Platyrrhines (platyrrhine means "flat nosed") the nose is broad and nostrils are set wide apart, well separated by a broad septum, and point sideways; in Catarrhines (catarrhine: "downward nosed") the nostrils are set close together, point forwards or downwards, and are separated by a very narrow septum.

Some Old World monkeys, particularly those that have

adopted a ground-living way of life, such as baboons, mangabeys, and the mandrills of the subfamily Cercopithecinae, appear to have re-adopted a long snout during their evolution. This structure, however, is not primarily olfactory in function but seems, rather, to be more closely related to the large size of the jaws and the prominence of the canine teeth; it should be considered a dental muzzle rather than an olfactory one.

Sensory reception and the brain. Among mammals in general, the olfactory system is the primary receptor for environmental information; consequently the brain of most mammals is dominated by the olfactory centres. In primates the sense of smell is considerably less important than the well-developed visual system and highly refined sense of touch. The primate brain is enlarged in the specific areas concerned with vision (occipital lobes) and touch (parietal lobes), thus taking a characteristic shape throughout the higher primates.

Touch. The skin of the primate hand is well adapted for tactile discrimination. Meissner's corpuscles, the principal receptors for touch in hairless skin, are best developed in the apes and man, but can be found in all primates except the tree shrew; structurally correlated with a high level of tactile sensitivity are certain anatomical features of the skin of the hands and feet, such as the absence of pads on the palms and soles, characteristic of mammals in general, and the presence of a finely ridged pattern of skin corrugations known as dermatoglyphics (the basis for fingerprints).

Eyes and vision. The evolutionary trend towards frontality of the eyes has not proceeded as far in the prosimian as in the anthropoid families. Tree shrews retain set of the eyes like that of the primitive ancestor, the central axis of each bony orbit being 140° apart. In lemurs this angle is considerably less, 60° – 70° , and in the anthropoid families of apes and monkeys the divergence has been reduced to 20° . It should be noted that the axes of the eyeballs (as distinct from the bony orbits) in apes and monkeys are, in fact, parallel.

Colour vision is of considerable advantage to arboreal animals living on fruits and insects. Species that have both rod and cone receptors in their retinas include all of the Anthropoidea with the exception of some of the nocturnal species (*Aotus*, the night monkey of South America, for example), and the family Loriscidae and some of the Lemuridae among the prosimians. Rods are the more primitive retinal cells and respond to low intensity light, while cones, which respond to high light intensity, have a greater resolving power, and thus provide for a finer type of visual discrimination. Cones, moreover, are adapted for colour discrimination.

Nervous integration. The elaboration of touch and vision supplements the senses of smell, hearing, and taste, providing the primate with a sensory armament of great range and flexibility. The primate central nervous system is sufficiently refined to deal with the elaborate bombardment of environmental information reaching it. Association areas provide the connections between the input and output centers of the brain—the motor and sensory cortex. Association areas are the memory banks where the memory of past experience is encoded in the infinitely complicated plexiform arrangement of the neurons, the

Colour vision

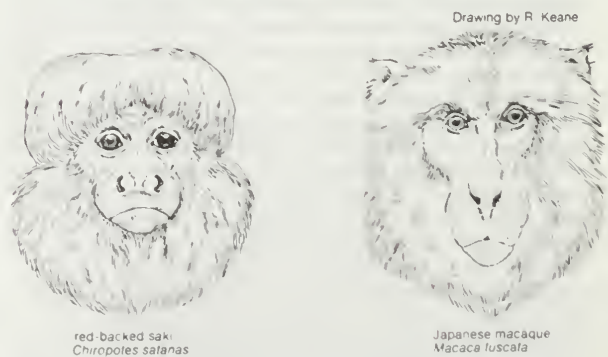


Figure 24. Head of New World monkey (saki) compared with head of Old World monkey (macaque).

brain cells, and their processes. All sensory impulses reaching the cortical centers of the central nervous system are routed through the association areas for conditioning, as it were, before reaching the effector side or output side where the appropriate response is initiated in the cells of the motor cortex. The more highly developed the association areas of the brain are, the more specific and appropriate is the behaviour and the more versatile is the animal in facing environmental demands.

The brain. The principal evolutionary trend in brain development has been towards elaboration. The neocortex of higher primates possesses highly developed associative functions, an aptitude for receiving, analyzing, and synthesizing the sensory input from visual, olfactory, auditory, gustatory, and tactile receptors and converting them into the appropriate motor responses.

The brain of Anthropoidea is larger, both absolutely and relatively, than that of prosimians. For instance, the weight of the simplest anthropoid brain, that of a marmoset, is three times greater than the brain weight of the bush baby, a prosimian of comparative size. This quantitative increase, in part, is attributable to the elaboration of the regions of the neocortex concerned with tactile and visual sensitivity; and in part to the elaboration of the intrinsic pathways connecting one part of the brain with another. The large brain of man is attributable not so much to an increased nerve cell content as to an increase in the size of the nerve cells and to a greater complexity of the connections linking one cell to another.

The external form of the anthropoid cerebral cortex is characterized by a complicated pattern of folds and fissures (sulci and gyri) in the brain surface. The fissural pattern is seen in its simplest form in the marmosets, but in the larger and more advanced members of the New World monkey infraorder, the capuchins (*Cebus*) for instance, the cerebrum is richly convoluted. Gyri and sulci are well marked in Old World monkeys and in the apes, the complexity of the pattern closely approximating the tortuous maze-like pattern seen in man.

Reproductive system. *Male and female genitalia.* The functions of the individual organs of reproductive systems are fairly uniform throughout the primates, but in spite of this physiological homology there is a remarkable degree of variation in minor detail of organs between groups, particularly in the external genitalia, which, by their variation, provide a morphological basis for the reproductive isolation of the species. There could be no more effective barrier to mating between different species than incompatibility of the male and female sex organs.

Among the characteristics of the primate order as listed by Mivart, the penis is described as "pendulous" and the testes as "scrotal." In contrast to other mammals, the primate penis is not attached to the abdominal wall but hangs free. The testes, with a few exceptions among the prosimians, in which they are withdrawn seasonally, lie permanently in the scrotal sac. In all primates except modern man (there is no evidence one way or the other regarding early representatives of the genus *Homo*), tarsiers, and some South American monkeys, the penis contains a small bone called the baculum, a typically mammalian character. The uterus of female primates shows all grades of transition between the two-horned (bicornuate) uterus, typical of most mammals, to the single-chambered (simplex) uterus of the higher primates and man.

Variations between primate taxa are demonstrated most strikingly by the glans penis, the scrotum, and perineum of the male, and the clitoris and the labial folds of the female vulva. The length and form of the clitoris, which when elongated mimics the penis, as in spider monkeys for instance, is a potent source of confusion in determining the sex of certain New World primates. The coloration of the male scrotum in forest-living primates, particularly of the genera *Cercopithecus* and *Mandrillus*, shows an infinite range of variations and provides a species-recognition signal of considerable effectiveness.

The external appearance of the genitalia undergoes seasonal variation in a number of primates. In the male, swellings of the testis and colour changes of the scrotum

occur, and, in the female, swelling and coloration of the vulva and perineal region herald ovulation, sometimes most obtrusively. Turgidity and excessive vascularity of the tissues of the perineum are probably characteristic of all mammals, but there are certain primate species in which this engorgement reaches monstrous proportions, notably baboons, mangabeys, some macaque species, and chimpanzees. Regions other than the primarily genital areas may also be affected by the hormones circulating at certain periods of the reproductive cycle. For instance, in the gelada baboon, the skin on the front of the female chest, which normally bears a string of caruncles resembling the beads of a necklace, becomes engorged and brightly coloured. A German zoologist, Wolfgang Wickler, has suggested that this is a form of sexual mimicry, the chest mimicking the perineal region. The observation that geladas spend many hours a day feeding in a sitting posture provides a feasible, Darwinian explanation of this curious physiological adaptation.

Placenta. The placenta, a characteristic of all Eutherian mammals, is a vascular structure that permits physiological interchange of blood and body fluids between the mother and the fetus and the breakdown products of the fetal metabolism; it also provides a two-way barrier preventing the passage of some, but not all, noxious substances and organisms such as bacteria and viruses from one individual to the other and is the source of hormones such as estrogens.

The placenta is a flat discoid-shaped "cake" in man, some of the other higher primates, and the tarsier. In many monkeys it is bidiscoidal, having two linked portions. The placenta is intimately attached on its outer surface to the endometrium, the lining of the uterus, by fingerlike processes (villi) that embed themselves in the endometrium where the vascular connections between the two circulations are achieved. The connection between fetal and maternal circulations appears as two distinct types among primates, a distinction which is believed to have had an important effect on the evolution of the order. In the first type (epitheliochorial), found in the prosimian suborder, several cellular layers separate the maternal and fetal bloodstreams, thus limiting the passage of molecules of serum proteins. In the second type, hemochorial placenta, the relationship is much more intimate, there being no fetal and maternal cellular layers separating the two circulations, and serum proteins can easily pass. A hemochorial type of placentation is found in tarsiers, apes, monkeys, and man.

The placenta is shed at birth in all primates, and, in all other than civilized man, is eaten by the mother.

EVOLUTION AND PALEONTOLOGY

Renewed interest in primate origins. Beginning in the 1950s there was a notable expression of interest in primate paleontology. Since then, hardly a year has passed without the announcement of some new major discovery. New sites have been opened up and old discoveries redescribed and reallocated. New techniques in geological dating, in palynology (the study of fossil pollen), in paleoclimatology and paleoecology and in the archeological interpretation of fossil sites, have helped to lift primate paleontology into the forefront of the life sciences and have aroused public interest to an unprecedented level. The popularity of all aspects of the evolution of man is reflected, for instance, in the spate of books covering this field published during this period.

The African continent has contributed the greatest share of significant early finds—Rusinga and Songhor in Kenya, Olduvai in Tanzania, and Sterkfontein, Kromdraai, Swartkrans, Makapan in South Africa are names with which every anthropology student is familiar. Later, the names Omo, Kanapoi, Lothagam, Napak, Moroto, Laetoli, and Hadar—all East African sites of hominid and pongid discoveries in the 1970s and '80s—became equally well known.

Elsewhere, pieces of this colossal worldwide jigsaw have been discovered in Europe, notably in Hungary, France, and Italy; in the Siwalik Hills of northwest India; in the ever prolific middle Eocene Bridger Beds of North Amer-

Mimicry of bottom by chest colours

Important recent sites of primate fossils

Relative versus absolute brain size

ica; and in Colombia and Bolivia. These latter discoveries were particularly welcome because of the poverty of fossils of New World monkeys.

While new discoveries have clarified the human story, older ones, that had only served to cloud it, have been repudiated. Piltdown man was shown unequivocally to be a fake in 1953; *Oreopithecus* of Monte Bamboli, Italy, has been relegated to a family remote from human ancestry, and Galley Hill Man in England has been shown to be a recent intrusion (a burial) into middle Pleistocene gravels. The questionable finds from the remoter geological period of the Eocene and Oligocene have also been re-examined, with the result that a number of confusing fossils have been dismissed. *Anagale*, long thought to be a fossil tree shrew, is now regarded as a small carnivore; others, such as the supposed tarsoid-anthropoid annectant form, *Parapithecus fraasi*, have been reallocated within the order Primates.

Progress in constructing the phylogeny of higher primates and man has been bedeviled by a number of controversies concerning taxonomy and nomenclature, and, as a result, internal disagreements have developed which mitigate against the sort of international cooperation that is essential for so vast a project as the construction of a human phylogeny. New-school and old-school taxonomists have come into conflict. But with the rapid advances in genetics, in the new concepts of biological species and in population anthropology, a fresh equilibrium is slowly being acquired, as the pendulum swings between the reactionary "lumping" of such taxa as genera and the traditional "splitting," in which every new discovery was provided with a new generic name. The most extreme example of taxonomic lumping in primate paleontology was in a 1965 revision by E.L. Simons and D.R. Pilbeam of the Miocene-Pliocene genus *Dryopithecus*; from the 28 genera recognized prior to 1965, the authors recommended the recognition of three genera only: *Dryopithecus*, *Ramapithecus*, and *Gigantopithecus*. All the rest were absorbed into the three genera, a result initially satisfying but not yet subjected to the test of time.

The primate fossil record. *Cretaceous.* The known range of the primate order was extended back from the middle Paleocene to the Late Cretaceous (about 75,000,000 years ago) by the discovery in Montana of one premolar and four molar teeth, representing two species of insectivore-like primates that were assigned in 1965 to a new genus, *Purgatorius*. This diagnosis is based on the characters of four molar and one premolar tooth and is not by any means universally accepted.

Paleocene. The best known early primate species, of which, *inter alia*, a complete skull and a partial postcranial skeleton are available, belongs to the genus *Plesiadapis* from the Paleocene (about 60,000,000 years ago) of Cenay, France, and Colorado. The skull shows a number of dental specializations including procumbent rodentlike incisors in the upper and lower jaw and the absence of antemolar teeth, but the molar teeth show the unequivocal primate affinities of this species. Rodent affinities have been suggested for *Plesiadapis*, but most authorities agree that it is a primate, however aberrant.

Eocene. The known Eocene fossil families include the Tarsiidae (tarsiers), the Adapidae, considered probable ancestor of lemurs and lorises, and the Omomyidae, possible ancestors of the monkeys and apes.

The family Adapidae contains two North American genera, *Notharctus* and *Smilodectes*, which are well represented in the fossil deposits of the Bridger Basin, Wyoming. *Adapis*, a closely related genus, was widely distributed in Europe, as were *Anchomomys* and *Pronycticebus*. *Notharctus* and *Smilodectes* are not thought to be antecedent to living lemurs, principally on zoogeographical grounds, since egress from the North American continent for tropical species was effectively cut off by ecological changes soon after the beginning of the Eocene (about 54,000,000 years ago). *Notharctus* was not unlike the modern lemurs in size and general appearance. On both morphological and zoogeographical grounds, *Adapis* and *Pronycticebus* or *Anchomomys* could have provided the stem from which the living lemurs and lorises evolved.



Figure 25: Eocene fossil lemur *Notharctus*.
Drawing by R. Keane

Representatives of the Omomyidae have been found in North America, in Europe, and in Asia. According to American paleontologist Elwyn L. Simons, the Omomyidae are probably the key family for tracing the phylogeny of the higher primates, the apes and monkeys. The family was common to North America and Eurasia, a fact which reinforces the presently held view that omomyids provided the stock from which New World monkeys and Old World monkeys, apes, and man arose. Their principal claim to antecedency (other than the circumstances of their zoogeography) is the dental formula: $\frac{2 \cdot 1 \cdot 3 \cdot 3}{2 \cdot 1 \cdot 3 \cdot 3} = 36$ teeth (18 pairs).

The Eocene Tarsiidae, represented by the European species *Necrolemur antiquus*, found in the Quercy deposits of France, are likely to contain the ancestor of the modern genus *Tarsius*. The tarsier is indeed a "living fossil" in the best sense of that overworked term. What is known of the skull and skeleton of *Necrolemur antiquus* indicates that it was almost identical with that of the living *Tarsius*.

Oligocene. Information on primate evolution during the Oligocene Epoch (about 38,000,000 to 26,000,000 years ago) rests principally on discoveries in three areas—Burma, Texas, and Egypt, and by far the most important is Egypt. From the Fayum region of the Sahara, from the upper levels of the Qatrani Formation, has come the first evidence of the emerging Anthropoidea. A number of different genera have been described from the Fayum, including *Apidium*, *Propliopithecus*, *Oligopithecus*, *Parapithecus*, *Aeolopithecus*, and *Aegyptopithecus*. The Fayum, indeed, seems to be the cradle of the Old World anthropoids. Simons has indicated that the possible ancestors of Old World monkeys (*Parapithecus*, *Apidium*), of gibbons (*Aeolopithecus*), of apes (*Aegyptopithecus*), and even of man (*Propliopithecus*) are all represented in the fauna. Unfortunately, apart from fragmentary scraps of limb bones, there are no postcranial bones to help. From the evidence provided by the possible forerunners and probable descendants of the Fayum primate fauna, however, it is reasonable to assume that quadrupedalism was becoming established as the typical locomotor pattern, and that vertical clinging and leaping, the characteristic gait of the Eocene forebears of the fauna, was probably retained by at least some of the genera represented at this site.

Of unusual interest is the recent discovery of the cranium of a North American omomyid called *Rooneyia*, it is of particular note in view of a belief that primates had disappeared from North America by the late Eocene times. *Rooneyia* is also of considerable interest in itself. The

Rooneyia,
an
important
transitional
form

Possible
ancestors
of the
lemurs

skull possessed a mixture of primitive (prosimian) and advanced (anthropoid) features, precisely the combination that might be anticipated in a transitional form between lower and higher primates.

Miocene. The Miocene Epoch is probably the most fruitful of all for paleoprimatology. During the 19 million years of its duration (beginning about 26,000,000 years ago), dramatic changes in geomorphology, climate, and vegetation took place. The Miocene was a period of volcanism and mountain building, during which the topography of the modern world was becoming established; of particular relevance to the story of primate evolution are the vegetational changes resulting from the orogenic ones. Grasses, known since the early Tertiary Era, flourished in the new conditions and in many areas that were previously forested. Grassland is known regionally by such names as savanna, llanos, prairies, etc. A new type of primate—the ground inhabitant—came into being during this period. The generalized nature of the bodily form of primates, combined with their specialized brain, made possible this critical step.

In the last few decades considerable additions to the knowledge of ape and human evolution have accrued from Miocene fossil beds in East Africa and Europe.

In Europe further discoveries of the gibbon-like genus *Pliopithecus* and the gorilla-chimp-like genus *Dryopithecus* have been made. *Pliopithecus* remains from Neudorf-ander-Mareh in Czechoslovakia have provided a remarkably complete picture of the habits of the Miocene apes, which, on this evidence, appear to have possessed bodily forms of a quadruped with morphological overtones of incipient brachiating specializations. Perhaps most interestingly of all, *Pliopithecus*, with its retention of numerous platyrrhine or New World monkey characters, emphasizes the phylogenetically conservative nature of anthropoid evolution.

In East Africa the excavations of the inshore islands and of the Kenyan shores of Lake Victoria by L.S.B. Leakey and a number of colleagues have illuminated man's knowledge of the evolution of the African great apes and have set up provocative questions regarding the origins of man. The *Proconsul* group (relegated to subgeneric status within *Dryopithecus* by the revision of Simons and Pilbeam) is known from three fossil species of ape, *Proconsul major*, *P. nyanzae*, and *P. africanus*. All three species are represented by a large number of jaw and facial fragments, but the only complete skull, discovered in 1948, belongs to *P. africanus*. Though complete, it was somewhat distorted by pressure of the surrounding rocks during fossilization. Subsequent reconstruction revealed a skull more monkey-like than apelike in its contours, an appearance which belied the unquestionably chimpanzee-like affinities of its teeth. The form of the skull and braincase of *P. africanus*, along with the forelimb skeleton, which is known in great detail for this species, indicates a body form that most closely resembles that of living monkeys. An American anthropologist, Sherwood L. Washburn, has described these creatures as "dental apes." The monkey-like characteristics of the *Dryopithecus* (*Proconsul*) group reflect the recency of the separation of the monkey and ape stocks of the Old World.

Proconsul africanus has been suggested as a possible common ancestor for apes and monkeys, but the general consensus is that *P. africanus* was already too specialized dentally to fit such a role. Be this as it may, there is no doubt that this Miocene fossil ape provides a very reasonable "model" for such an ancestor.

In East Africa two species of fossil gibbons, *Limnopithecus*, have also been discovered, closely related, if not generically identical, to the fossil "gibbon" of European deposits termed *Pliopithecus*. The taxonomic status of these fossils as "gibbons" is somewhat uncertain.

It is towards the end of this turbulent epoch that the first probable ancestral hominids appear in the fossil record. *Ramapithecus punjabicus* was discovered in 1934 by an American paleontologist, G.E. Lewis, in the Siwalik Hill Miocene deposits in northwest India. Lewis tentatively expressed the view that *Ramapithecus* was a hominid, but the scientific climate of the period was not propitious for claims that seemed to imply such remote antecedents for

man; as a result Lewis' perceptiveness was largely ignored. Recently, Simons and Pilbeam have resurrected Lewis' specimen and have demonstrated its hominid characteristics. At Fort Ternan in East Africa, meanwhile, Leakey found what at present appears to be another specimen of *Ramapithecus*, but of a different species and called by him *Kenyapithecus wickeri*.

It was also from the late Miocene that the first unequivocally ground-living monkey came to light. Attributed to the living subfamily Cercopithecinae of Old World monkeys, these specimens of *Mesopithecus pentelici* demonstrate the transition between the ancestral monkey type (today represented by the leaf-eating monkeys or Colobinae of Africa, India, Ceylon, Malaya, Burma, and Southeast Asia) and the newly evolved type of ground-living monkey epitomized by *Macaca* or *Papio*. In its known features of skull and postcranial skeleton, *Mesopithecus* is broadly equivalent to the macaque monkeys.

Pliocene. The Pliocene Epoch (from about 7,000,000 to 2,500,000 years ago) was very similar to the present in terms of its geomorphology and climate. Discounting the effects of the recent influence of man on the distribution of forest and savanna in the tropics, the face of the land cannot have differed much from the aspect it presents today. Thus, one would expect that during the Pliocene (bearing in mind the effectiveness of environmental selection) essentially modern forms of primates would have made their appearance, but perhaps the most famous of the Pliocene fossils was far from "modern." *Oreopithecus* was first discovered, fragmentarily, in brown coal deposits in Italy in the middle of the 19th century. The most recent discovery of a complete, though paper-thin, skeleton (so grossly has it been compressed) in Tuscany has demonstrated the power of the environment acting upon the genetic structure, in once again reproducing the trade marks of two of the "inevitable" trends of primate evolution, brachiation, and bipedalism. *Oreopithecus* possessed a number of dental and bony characters that are typically hominid; it also possessed limb characteristics that are typically pongid. The canines were relatively short, the face was abbreviated, and the pelvis was broad. The foot showed characteristics associated with bipedal walking as did the vertebral column (all hominid characters), yet the arms were long and the fingers long and curved. The limb proportions are those of a brachiator (pongid characters). Finally, a number of characters of this confusing hominoid typically reflect those found in Old World monkeys. The present consensus of attitude towards this composite primate is that it is best segregated within the superfamily Hominoidea but in a family of its own, the Oreopithecidae. That the end of the *Oreopithecus* story has not yet been heard is certain; in the meantime the Tuscany fossil stands as a memorial to the primate order, the order which, of all of the class Mammalia, is most subject to environmental adaptation by virtue of its behavioral plasticity.

Few fossils referable to the ape family are known during the Pliocene, and monkey families are scarcely better known. *Libypithecus* and *Dolichopithecus*, both monkeys, were probably ancestral cercopithecines that still retained some features of their colobine ancestry; but neither genus can be placed in a precise ancestral relationship with modern members of this subfamily.

Pleistocene. The Pleistocene, which on faunal evidence is estimated to have started 2,000,000 to 3,000,000 years ago, is the epoch of hominid expansion (see EVOLUTION, HUMAN: *Classes of Hominidae*). Knowledge of nonhuman primates is surprisingly sketchy. No ape fossils are known until relatively recent times and monkeys have been identified in only a few regions in Africa and even fewer in Asia; e.g., *Cercopithecoides* (a colobine-like form), *Gorgopithecus*, *Dinopithecus*, from South African deposits (forms of the living genus *Papio*); and *Simopithecus* (a giant, ancestral forerunner, according to one authority, of the present-day genus, *Theropithecus*, the gelada) from Olduvai Gorge and South Africa. It is possible that the *Papio-Theropithecus* divergence can be pushed well back into the Pliocene.

One genus of the Pleistocene that is neither ape nor

First
ground-
living
monkey

The special case of *Gigantopithecus*

monkey in the sense that these taxa are interpreted today is *Gigantopithecus*. The romantic story of the discovery of the gargantuan molar teeth of *Gigantopithecus blacki* by von Koenigswald, the Dutch paleontologist, in a Chinese pharmacy has often been told. The boldness of the move that erected a new genus on such apparently slender grounds has been amply justified by the subsequent discovery of several massive jaws from Kwangsi in South China and attributed to this genus. The most recent discovery is of an enormous jaw in the Dhok Pathan deposits of the Siwālik Hills of India. Possibly dating back to the middle Pliocene, this discovery provides a respectably long period of existence for this aberrant giant-toothed hominoid genus. Clearly *Gigantopithecus bilaspurensis* was a member of the family Pongidae with divergent dental specializations that were possibly adaptive for foraging in grassland where tree products were unavailable and ground products available but hard to get. *Gigantopithecus* was a fossil ape that fell by the wayside on its way up to hominid status, a victim, possibly, of overspecialization.

CLASSIFICATION

Distinguishing taxonomic characters. The identifying characters of the primate order have already been discussed above. The point has been made that no single morphological characteristic distinguishes this group exclusively from other mammalian taxa. For instance, while it is possible to make the generalization that primates bear flat nails on the ends of their fingers and toes, it is not strictly accurate, for species such as the tree shrews, the marmosets, and the tamarins have either overt claws (tree shrews) or modified claws (marmosets and tamarins). One can also aver that primates characteristically possess two pectoral mammae, but it must be pointed out that the aye-aye (*Daubentonia madagascariensis*) has its nipples sited in the inguinal region (lower lateral region of the belly) and that mouse lemurs and tree shrews have multiple sets. Generalizations regarding opposability of the thumb, a "classic" primate character, are particularly vulnerable since some marsupials possess this trait and some primates do not. Examples of the generalized nature of primate specializations are too numerous to mention here. At first sight, big toe opposability provides a more promising diagnostic feature but, as with other features, there are exceptions to the rule, man for instance; moreover, it occurs in nonprimates.

As F. Wood Jones said in his book *Man's Place Among the Mammals*: "There is no single [characteristic] which constitutes a peculiarity of the Primates: for a primate animal may only be diagnosed by possessing an aggregate of them all." It was with this precept in mind that Sir Wilfrid Le Gros Clark introduced his list of evolutionary trends of primate evolution in 1959, which is amplified below:

1. The preservation of a generalized structure of the limbs with a primitive pentadactyly [*i.e.*, containing five digits] and the retention of certain elements of the limb skeleton (such as the clavicle) which tend to be reduced or to disappear in some groups of mammals.
2. An enhancement of the free mobility of the digits, especially the thumb and big toe (which are used for grasping purposes).
3. The replacement of sharp compressed claws by flattened nails, associated with the development of highly sensitive tactile pads on the digits.
4. The progressive abbreviation of the snout or muzzle.
5. The elaboration and perfection of visual apparatus with the development to varying degrees of binocular vision.
6. Reduction of the apparatus of smell.
7. The loss of certain elements of the primitive mammalian dentition, and the preservation of a simple cusp pattern of the molar teeth.
8. Progressive expansion and elaboration of the brain, affecting predominantly the cerebral cortex and its dependencies.
9. Progressive and increasingly efficient development of those gestational processes concerned with the nourishment of the foetus before birth. (From W.E. Le Gros Clark, *The Antecedents of Man*; © 1959 Edinburgh University Press.) Two further evolutionary trends may be added to the list:
10. Progressive development of truncal uprightness leading to a facultative (opportunistic) bipedalism.
11. Prolongation of postnatal life periods.

From J.R. and P.H. Napier, *A Handbook of Living Primates* (1967).

A British zoologist, R.D. Martin, has criticized many of these evolutionary trends as irrelevant as determinants of primate origins. In his view, many of them are either symplesiomorph characters (primate characters inherited from an ancestral mammal and therefore nonspecific for primates) or convergent characters (characters that have been acquired by many groups of primates in response to the environmental effects of natural selection). The disputed trends of primates, nonetheless, comprise characters that are of considerable significance for primate classification. Martin's point is that such characters are of little use in distinguishing the order Primates from other placental orders. Martin argues that the inadequate nature of the current definitions of primates is the principal reason why, when so much behavioral and structural evidence points to the opposite conclusion, the tree shrews (Tupaiaidae) are presently included in the order. This may well be so. But until the true synapomorph characters (characters of independent acquisition) are identified by extensive comparative studies of all of the relevant characters of primates and nonprimates, there is little justification for arbitrary disenfranchisement of the tree shrews. "Possession is nine-tenths of the law," a principle that, as in this instance, supplies the highly desirable element of stability to zoological nomenclature.

Annotated classification. It is apparent that the taxonomy of primates is in a very fluid condition; but this state of affairs is wholly preferable to the virtual stagnation that existed during the first half of the present century. In view, therefore, of the likelihood of further major revisions taking place during the next decade that may well reveal the fallacies of certain current concepts, the informational content of this section is best served by a conservative, rather than a radical, approach to primate classification. Thus, the most widely accepted current classification is here followed, while fully admitting that it is not wholly adequate. In order to avoid some of the descriptive complications, derived from an imperfect classification, the following scheme outlines the distinguishing taxonomic characters of order, suborders, and families; intermediate taxa, such as infraorders and superfamilies, are omitted.

ORDER PRIMATES

Distinguished by the possession of a hallux (first toe) and a pollex (thumb), one or both opposable; eyes frontally placed on the face, orbit ringed by bone and, in higher members of the group, closed posterolaterally by bone. Brain relatively large. Most of the fingers and toes bear nails, not claws. Number of teeth 18 to 36; canines prominent but dentition as a whole not highly specialized. Diet omnivorous; body weight from about 100 grams (a few ounces) to 180 kilograms (400 pounds), more or less. Twelve Recent families, 55 genera and 197 species; 8 extinct families. Extinct groups are indicated by daggers (†).

†Suborder Plesiadapoidea

†*Families Phenacolemuridae, Carpolestidae, and Plesiadapidae*
Middle Paleocene to lower Eocene; North America and (except Phenacolemuridae) Europe; collectively, more than a dozen genera.

Suborder Prosimii

Usually quadrupedal and plantigrade, but some species show vertical clinging in trees and saltation on the ground; all have long tails (except *Indri*). Very small to medium in size; substantially insectivorous, but fruit- and leaf-eating common; activity rhythm nocturnal or crepuscular. Suborder characterized by possession of long or longish snouts, naked rhinarium, and naked tethered upper lip (except *Tarsius*); facial vibrissae (whiskers) present. Specialized sublingual fold, often serrated, below tongue. Eyes large, fully frontally facing in *Tarsius*, but somewhat divergent laterally in others; bony orbits not fully closed posterolaterally. Teeth characterized by special arrangement of lower incisors and canines to form a dental "comb," except in *Tarsius* (structure absent) and in *Tupaia* (canines take no part in comb); upper incisors usually small and peglike. Hallux and pollex more or less opposable, tips of all digits bear flat or flattish nails (except in *Tupaia*), second digit of hindlimb bears toilet (grooming) claw. Placenta epitheliochorial and non-deciduate (except *Tarsius*). One or more pairs of nipples. Single births are the rule.

Family Tupaiaidae (tree shrews)

Small animals with elongated body and long bushy tail (sca-

Features and exceptions

ly with terminal tuft in *Ptilocercus*). Short limbs, with fore longer than hind. Long pointed snout. Relatively large brain (compared with nonprimates) with diminished olfactory representation; turbinal bones of nasal region of skull not well developed. Tympanic ring within middle ear cavity. Primitive dental pattern, with 3 pairs of lower incisors, protuberant and arranged as modified "dental comb"; cheek teeth primitive in form. Digits of hands and feet bear claws; pollex is incipiently divergent, but nonopposable; hallux, nonopposable. Both arboreally and terrestrially adapted species occur within the family. Head and body length 10–24 cm (3.9–9.4 in.), tail 12–14 cm (4.7–5.5 in.); 5 genera, 17 species; India and Southeast Asia and Philippines.

†*Family Adapidae* (Notharctidae)

Lower to upper Eocene; Europe, Asia, and North America; about 10 genera.

Family Lemuridae (lemurs)

Size very small to medium; tail longer than head and body length combined; hindlimbs considerably longer than forelimbs; tail more developed in subfamily Lemurinae than in the other subfamily Cheirogalinae. Coat colours highly variable; one species, *Lemur macaco*, shows sexual dichromatism (males black, females rufous). Tympanic bullae large, tympanic ring contained within middle ear cavity (contrasting with Lorisidae). Cutaneous glands on forearm and upper arm serving marking functions in *L. catta* and *Hapalemur griseus*; absent from other species. Dentition typical of Prosimii, including "dental comb" and reduced upper incisors (absent in *Lepilemur*). Diurnal (*Lemur*, *Hapalemur*) or nocturnal. Births seasonal; twins common in *Microcebus*. Head and body length 13–46 cm (5.1–18.1 in.), tail 17–56 cm (6.7–22 in.); 6 Recent genera; about 14 species (some possibly extinct); restricted to Madagascar.

Family Indridae (indri, sifaka, and avahi)

Arboreal animals with exceptionally long limbs, associated with specialized locomotor habit of vertical clinging and leaping. Skull with somewhat downwardly facing foramen magnum, short facial skeleton and very large auditory bullae. Dentition atypical of Prosimii; lower canines absent, first premolar somewhat caniniform, 2 premolars only. Anatomy of stomach and cecum, particularly, specialized in association with a bulky vegetarian diet. Activity rhythm nocturnal (*Avali*) or diurnal (*Propithecus* and *Indri*). Fossils from Pleistocene. Three living genera, 4 species; Madagascar.

Family Daubentoniidae (aye-aye)

Arboreal, nocturnal family comprising single species *Daubentonia madagascariensis*. Hand specialized by markedly elongated and attenuated third digit. Dentition anomalous: central incisors reduced to single large pair having long roots and a rodent-like function; canines absent, premolars absent in lower jaw and reduced in upper. Symphysis between right and left halves of mandible is fibrous, a condition unique among the primates. Single inguinal pair of nipples. Madagascar.

Family Lorisidae (galagos, lorises, angwantibo, and potto)

Arboreal and nocturnal. Two major locomotor types comprise family: vertical clingers and leapers (subfamily Galaginae) and slow quadrupedal climbers (Lorisinae), the former possessing excessively long hindlimbs, and the latter subequal limbs. Tails short or absent in Lorisinae and long in Galaginae. Second digit of hand reduced, extremely so in Lorisinae (particularly *Perodicticus* and *Arctocebus*) and first digit of hand and foot disproportionately large. Skull globular in form with little protrusion of facial region. Tympanic bullae large, short external auditory canal with tympanic ring sited at junction of middle and external ear. Orbits large, with little lateral divergence. Dentition typical of suborder. Omnivorous (with large insectivorous element). Family comprises 5 Recent genera with about 12 species; Africa, Asia (Lorisinae); Africa (Galaginae).

†*Family Anaptomorphidae*

Upper Paleocene to middle Eocene; North America and Europe; about 10 genera.

†*Family Omomyidae*

Upper Paleocene to lower Oligocene; North America, Europe, and Asia; about 20 genera.

Family Tarsiidae (tarsiers)

Arboreal and nocturnal. Hindlimbs almost twice the length of forelimbs. Hands with 5 digits, fingers surmounted by circular digital pads and backed by tiny triangular nails; thumb opposable, but in a fashion distinct from other Prosimii. Head of femur cylindrical (not spherical as in all other Prosimii except *Galago*), and tibia and fibula fused in lower third (unique to *Tarsius*). Skull characterized by enormous, frontally directed orbits, incompletely closed posterolaterally; nasal space much compressed. Foramen magnum (opening for spinal cord) di-

rected downwards. Short external auditory canal; middle ear with prominent bulla and tympanic ring, sited as in Anthropoidea. Dentition unlike typical prosimian pattern, lacking dental comb; central lower incisors absent; molars sharp-cusped tuberculo-sectorial. Placentation is hemochorial, discoidal, deciduate, as in Anthropoidea. Single young usually produced. One genus, 3 species; Southeast Asia.

†*Family Microsyopidae*

Lower to upper Eocene; North America and Europe; 4 genera.

Suborder Anthropoidea

Quadrupedal, brachiating or bipedal in gait. Diurnal in habit (except night monkey, *Aotus*). Omnivorous in diet. Tail present or absent, varying in length from species to species, prehensile in some. Size ranges from very small (*Saguinus* and *Callithrix*, the tamarins and marmosets) to very large (*Gorilla* and *Homo*.) Faces show varying degrees of protrusion, most marked in the baboons and mandrills and least developed in man. Neurocranium relatively large and with or without sagittal or nuchal crests or both; foramen magnum directed downwards and backwards. Orbits frontally facing and fully closed posterolaterally; nasal cavity narrow. Dentition comprising 30–36 teeth, canines are well developed except in man. Hands and feet prehensile; thumb semi-opposable or fully opposable (except in certain platyrrhine species). Hallux widely abductable (capable of being rotated inward) and capable of prehension in all species except man; flat nails on fingers and toes (except Callitrichidae); specialized toilet claws absent. Two pectoral mammae only. Placenta hemichorial, discoidal, and deciduate. Lengthy gestation period and extended life periods (compared with Prosimii) characterize the suborder. Suborder contains two geographical groups: New World monkeys and Old World monkeys; 6 families and 35 genera; both the tropical and subtropical zones of New and Old World.

Family Callitrichidae (tamarins and marmosets)

Small animals, diurnal, arboreal, and mainly insectivorous. Quadrupedal gait, hindlimbs being 25 percent longer than forelimbs. Thumb nonopposable, but hallux fully abductable; all digits bear modified claws (tegulae) except hallux, which has a flat nail. External nostrils face forwards and sideways (platyrrhine condition); nasal septum broad. Coat colour highly variable between species including wholly black and wholly white forms; one species (*Leontopithecus rosalia*) golden; distinguishing external features include a variety of ear tufts and a diversity of mustaches. Skull long and ovoid in shape with bulging occiput; tympanic bulla present; external auditory canal absent. Dentition consists of 32 teeth (36 in *Callimico*) with 3 premolars and only 2 molars. Tamarins are distinguished from marmosets by long lower canines which project beyond incisors; in marmosets canines and incisors are the same length. Head and body length 13–37 cm (5.1–14.6 in.), tail 20–39 cm (7.9–15.4 in.); weight about 100–700 gm (3.5–24.7 oz). Five genera, 15 species; South America.

Family Cebidae (New World monkeys)

Larger in size than Callitrichidae, with ears more or less naked externally. All digits bear flattened or curved nails; thumb imperfectly opposable. Gait quadrupedal. Prehensile tail present in *Ateles*, *Alouatta*, *Brachyteles*, and *Lagothrix*. Bony ear formation as in Callitrichidae. Nose typically platyrrhine. Dentition of 36 teeth; molars quadrangular and quadricuspid, and canines moderately developed in both sexes. Head and body length about 24–70 cm (9.4–27.6 in.); tail usually of same magnitude, weight about 600 to 10,000 gm (1.3 to 22 lb); 11 genera, 29 species; southern Mexico to Brazil.

†*Family Parapithecidae*

Lower Oligocene; North Africa; 2 genera.

Family Cercopitheidae (Old World monkeys)

Medium to large in size. Arboreal or terrestrial. Characterized by naked callosities covering expanded ischial bone of pelvis. Quadrupedal, plantigrade, or digitigrade in gait; hindlimbs longer than forelimbs in most species, but in ground-adapted forms limbs are nearly equal. Thumb present in all genera (except *Colobus*) and functionally opposable. Cheek pouches present in subfamily Cercopithecinae and concerned with temporary storage of food; in subfamily Colobinae pouches are absent but stomach is sacculated in association with a predominantly leaf-eating diet. Tail length variable from very long (Colobinae) to very short or absent (some Cercopithecinae). Muzzle more or less prominent. Noses with external nasal opening facing forwards and downwards or directly downwards ("catarrhine" condition); internasal septum narrow (as contrasted with Callitrichidae and Cebidae). Skull of anthropoid type but without auditory bulla; external auditory canal present. Dentition of 32 teeth; canines large in males, moderate in females; first lower premolar sectorial, molars quadricuspid with transverse linking crests (bilophodonty). Head and body length about 32–

110 cm (12.6–43.3 in.); tail, when present, 2–100 cm (0.8–39.4 in.); weight about 0.7–30 kg (1.5–66 lb); 13 genera, 72 species; Africa, India, China, Japan, and South East Asia.

Family Hylobatidae (gibbons and siamang)

Demonstrate highly specialized form of locomotion called brachiation. Arboreal, tailless catarrhines with forelimbs twice length of hindlimbs. Hands exceptionally long and prehensile; deep cleft between first and second digits; thumb fully opposable by means of unique mechanism at carpo-metacarpal (basal thumb) joint, associated with a specialized musculature. Ischial callosities as in Cercopithecidae but developing late in life. Well developed laryngeal (vocal) sac in *Symphalangus*. Skull of typical anthropoid form, orbits large, forwardly directed and having a well-butressed lateral margin; strong brow ridge; sagittal and nuchal crests rare; mastoid processes absent. Long sabre-like canines equally developed in both sexes; molar cusp pattern similar to that of the Pongidae and Hominidae. Chromosome diploid number: 44. Head and body length about 40–64 cm (15.7–25.2 in.); weight 4–13 kg (8.8–28.6 lb); 2 genera; *Hylobates* (6 species) and *Symphalangus* (1 species); Southeast Asia.

Family Pongidae (great apes: gorilla, chimpanzee, and orangutan)

Arboreal or ground living. Locomotion consists of brachiation and a modified form of quadrupedalism, called knuckle walking. Characterized by total lack of tail, relatively long forelimbs, elongated hands, and short but fully opposable thumbs. Brain relatively large compared with body size and showing well marked cerebral fissuration. Permanent dentition 32 teeth: incisors large and spatulate, particularly in upper jaw; canines large and conical, particularly in males; well-marked diastema present between upper lateral incisors and canines. Molars large; quadricuspid in upper jaw and quinquecuspid in lower; cusps arranged in so-called “*Dryopithecus*-pattern”; enamel much wrinkled in *Pongo*; third molar relatively small in *Pan*. Skull bears sagittal and nuchal crests in adult males (but rare in *Pan*). Ischial callosities usually absent (cf. Cercopithecidae, Hylobatidae). Thorax widest transversely (cf. Cercopithecidae); pelvis expanded in iliac region; lumbar vertebrae reduced in number (cf. Cercopithecidae). Life periods extended compared with New and Old World monkeys. Chromosome diploid number: 48. Family includes three genera: *Gorilla*, *Pan*, *Pongo*. *Distribution*: Africa and Southeast Asia.

Family Hominidae (man and fossil relatives)

Ground-living, hairless, omnivorous and bipedal. Hindlimbs longer than forelimbs (contrasting with Pongidae and Hylobatidae). Thumb fully opposable and hand capable of both precision grip and power grip; big toe aligned with other toes, foot nonprehensile. Brain relatively large and highly fissured. Neurocranium expanded particularly in parietal and frontal regions; Mastoid process of periotic bone present from birth. Facial skeleton nonprotrusive and strongly flexed relative to cranial base. Dentition of 32 teeth: canines small and blunt, lying flush with occlusal plane; premolars bicuspid; molars small, low, bearing rounded cusps; last molar in upper and lower jaw frequently unerupted. Postcranial skeleton shows adaptations to upright posture particularly in vertebral column, pelvis and lower limb. Chromosome diploid number: 46. Life periods (including growth period and life expectancy) extended as compared with Pongidae. Head and body length about 80 to more than 130 cm (31.5–51.2+ in.); total length (height) about 130 to over 220 cm (51.2–86.6 in.); weight from about 40 to more than 120 kg (88–264 lb); one Recent species, *Homo sapiens*; distribution world-wide.

Critical appraisal. The classification of the primates by George Gaylord Simpson is the most widely used list today. Since it was published in 1945 a number of revisions at family level and below have been incorporated; the higher categories have remained unchanged, although they have been the subject of considerable criticism and doubt. W.C. Osman Hill's classification, which, with modifications, follows R.I. Pocock's list published in 1918, has not enjoyed a wide popularity, in spite of the success of “Primates,” Hill's monumental treatise on the order. Hill excludes the tree shrews from the primates altogether; he also removes the tarsioids from among the lemurs and lorises and places them with the monkeys and apes. A compromise view between the classifications of Simpson and Hill has been proposed by American paleontologist Alfred S. Romer. Romer also omits the tree shrews from the order and gives the tarsioids equal status with lemurs and lorises, Old World monkeys, New World monkeys, and apes in a separate subordinal group. It seems likely

that before very long Romer's classification, or one closely similar to it, will become generally accepted. Reference to the annotated classification given above (based on Simpson's) will indicate how ill-fitted is the living genus *Tarsius* to share subordinal status with the prosimians.

Several revisions of lower taxa (genus and species) have been proposed in recent years which are proving generally acceptable. It has been suggested that *Cynopithecus niger*, the Celebes black “ape,” be included within the synonymy of *Macaca*, in which case this species therefore becomes *Macaca nigra*. Also widely accepted is inclusion of *Comopithecus* (the hamadryas baboon) in the synonymy of *Papio* (the common baboon), but opinions differ as to whether the hamadryas and the common baboon are not better considered a single species in view of the good evidence of hybridization in the wild. Synonymy of the genus *Mandrillus* with *Papio* has also been suggested, largely on the evidence of molecular data from blood protein studies, but has not found much support. Less controversial is the conclusion, based on field evidence, that the three generally accepted species of savanna monkeys (*Cercopithecus aethiops*, *C. sabaues* and *C. pygerythrus*) actually comprise a single polytypic species, *C. aethiops*.

Less universally acclaimed is the suggestion by Simpson that the structural and behavioral differences between the chimpanzee and the gorilla are insufficient to justify continued generic separation. According to this view, which has some support, *Pan*, having nomenclatorial priority, should include two species, *Pan troglodytes* (chimpanzee) and *Pan gorilla* (gorilla). At a symposium on the Old World monkeys held in 1969, the synonymy of *Nasalis* (the proboscis monkey) with *Simias* (the Pagai Island langur) and *Rhinopithecus* with *Pygathrix* was recommended but, as yet, is by no means universally accepted. Field, museum, and laboratory work has led to numerous changes in classification of primates by subspecific level that need not concern us here. (J.R.N.)

Edentata (armadillos, sloths, anteaters)

Edentata, an order of mammals, includes 31 living species distributed among the armadillos, true anteaters, and tree sloths, as well as eight extinct families of ground sloths and armadillo-like animals. The living families and six of the extinct families comprise the suborder Xenarthra. A second suborder, Palaeonodonta, consists of two extinct families. The entire evolutionary history of the edentates is restricted to the Western Hemisphere and the majority of the living species occur today in South America.

GENERAL FEATURES

Edentata means lacking teeth, though in reality only the true anteaters are toothless. The majority of edentates have simple, peglike cheek teeth that lack enamel; canine-like teeth do occur in some forms. Certain armadillos may have as many as 100 teeth. Edentates possess specialized traits, such as reduced dentition, a long sticky tongue, powerful, clawed, forefeet, associated with their insect diets. They also have primitive traits, such as the possession of five toes on the hindfeet, a simple uterus, and small, uncomplicated brain, which place them close to the primitive stock that gave rise to the infraclass Eutheria, which includes all placental mammals.

A persistent myth holds that any kind of armadillo, in defense against enemies, can roll itself into an impregnable armour-plated ball. In fact, only the three-banded armadillo (*Tolypeutes tricinctus*) is able to perform this feat. Most armadillos contort themselves into a reasonable facsimile of a ball, but some only slightly increase the convexity of the dorsal shell (carapace) while flattening against the ground. Individuals of some species flee to burrows or attempt to burrow right on the spot but may freeze and feign death when flight and rapid excavation are somehow prevented. Misconceptions about other edentates also exist. The toothless, weak-jawed anteaters seem to be sucking in ants and termites through a small tube, when, in fact, it is the sticky product of the greatly enlarged submaxillary glands covering the lengthy extensible tongues that facilitates the capture of these insects. The

Problem of the tarsiers

Primitive traits

slow-moving arboreal sloths, typically depicted as moving upside down along a branch, really spend only about 10 percent of their time moving at all.

The armadillos and anteaters are important ecologically, helping to control many kinds of abundant insects, especially ants and termites. Armadillos feed on some carrion and thus play a small part in cleaning up the environment. Although some armadillos may become agricultural pests, and the armour of several species is made into purses and other curios, edentates are generally of little economic importance. Generally, the greatest interest in edentates lies in their bizarre appearance.

Size range

Living edentates range in size from the tiny pink fairy armadillo or "lesser pichiciego" (*Chlamyphorus truncatus*) of Argentina, measuring about 16 centimetres (6.3 inches) in length and weighing little more than 100 grams (3.5 ounces), and the slightly larger two-toed or silky anteater (*Cyclopes didactylus*), just over 37 centimetres (14.6 inches) long and weighing about 325 grams (11.5 ounces), to the 60-kilogram (132-pound) giant armadillo (*Priodontes giganteus*), nearly 1.5 metres (five feet) long, and the giant anteater (*Myrmecophaga tridactyla*), which weighs up to 25 kilograms (55 pounds) and may be over two metres (6.6 feet) in length. Among extinct forms, an Oligocene armadillo, *Prozaedyus proximus*, had a skull under eight centimetres (3.1 inches) in length and was only somewhat larger than the smallest of the living armadillos. The Pleistocene ground sloth, *Megatherium americanum*, was six metres (20 feet) long and was larger than a modern elephant; a Pleistocene glyptodon (*Doedicurus clavicaudatus*) was more than four metres (13 feet) long and 1.5 metres (five feet) high. Both flourished in South America about 500,000 years ago.

Anteaters, tree sloths, and armadillos are strikingly different in appearance, but within each of these families, there is only a modest diversity of structure. The toothless anteaters vary from strictly arboreal to terrestrial forms that are characterized by progressively more elongate tubular mouths. The arboreal species have only two (*Cyclopes*) or three (*Tamandua*) clawed toes on the forefeet, whereas the terrestrial giant anteater has four. In addition, the arboreal species have prehensile tails. At a casual glance, the two genera of herbivorous (plant-eating) tree sloths (*Bradypus*, *Choloepus*) differ from each other in general appearance only in the number of toes on the forefeet (two in *Choloepus* and three in *Bradypus*) and in pelage coloration. The armadillos (Dasypodidae) are distinguishable on the basis of relative size and on the distribution, flexibility, and pattern of the bony shields that form their body armour.

The anteaters, essentially a tropical family, occur from northeastern Argentina and northern Uruguay to the extreme northern part of the state of San Luis Potosí, Mexico. Tree sloths are restricted to forests from northern Rio Grande do Sul, Brazil, to southeastern Honduras. Armadillos are found from the northern part of the province of Santa Cruz, Argentina, to south-central Kansas. Most living edentates are fairly common in their preferred habitats, although several species have disappeared from parts of their range largely because of habitat destruction and excessive hunting by man. These endangered species include the giant anteater, the three-toed sloth (*Bradypus torquatus*), the three-banded armadillo, the giant armadillo, the pink fairy armadillo, and Burmeister's armadillo (*Burmeisteria retusa*).

NATURAL HISTORY

Behaviour. Anteaters eat termites, ants, and occasionally other insects. They may forage in trees (two-toed), on the ground (giant *Myrmecophaga*), or in both situations, as does the lesser anteater (*Tamandua tetradactyla*). Their prey is usually located by smell, the sense of which is better developed in these animals than are those of eyesight and hearing. With the long, powerful foreclaws, the anteaters rip open anthills and termite mounds, into which they then insert their long sticky tongues. The tongue, covered with prey and debris, is then drawn into the tubular mouth. When frightened, the giant anteater usually flees at a slow clumsy gallop, but if forced it will fight, using

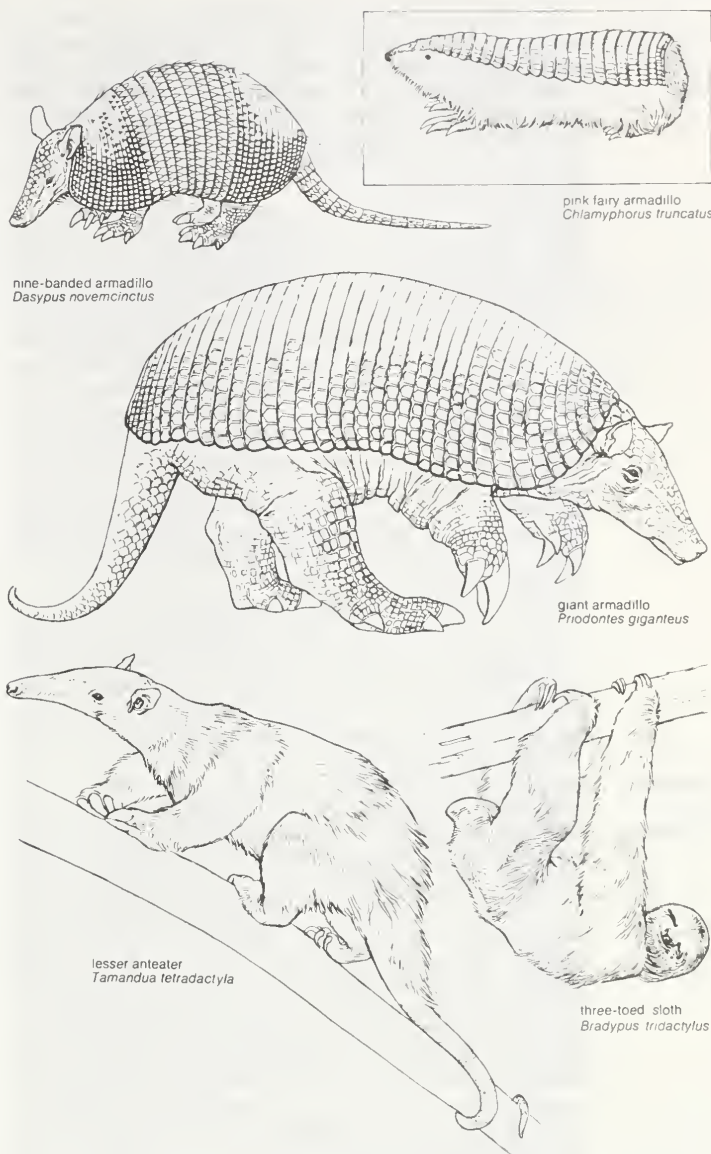


Figure 26: Edentate body plans.

Drawing by J.C. Barbers

its claws to rip the enemy. *Cyclopes* and *Tamandua* also use their foreclaws for defense, anchoring themselves to a branch with their prehensile tails. When alarmed, these two emit shrill calls.

Anteaters travel solitarily or as mother and young and do not seem to have permanent resting places. The giant anteater chooses a secluded spot such as a hollow log or the abandoned burrow of some other animal and then curls up to rest, with its long bushy tail covering its head and body. The arboreal anteaters find shelter in trees. The two-toed and lesser anteaters are mainly nocturnal, but the giant anteater is both nocturnal and diurnal. In areas remote from civilization, the giant anteater is most active during the day, whereas in heavily populated regions it is primarily nocturnal.

Tree sloths are strictly herbivorous, but the absence of enamel on their teeth suggests that the group perhaps was once insectivorous. The lack of enamel could pose serious problems relating to tooth wear for animals eating such abrasive food as vegetation. The difficulty, however, is alleviated by continuous growth of the teeth, which have an outer covering of hard cement. As is typical of herbivores, tree sloths have large multichambered stomachs in which the slowly digesting plant material steeps. Constantly filled, the stomach and its contents constitute about 30 percent of the sloth's body weight. A meal may take up to a week to digest. Every two to eight days, the sloth descends to the foot of its home tree and passes a large amount of excrement.

Food habits of sloths

The diet of the two-toed sloth consists of a variety of leaves, stems, and fruits. The three-toed sloth, on the other hand, eats only the leaves of the cecropia tree. The sloth locates its food by touch and smell and then hooks a branch with its curved claws and pulls it to the mouth. Although colour vision is present, the eyesight and hearing of sloths are not very well developed, and they orient themselves mainly by touch.

In the wild, the hair of the sloth's body is inhabited by algae, which flourish in the jungle humidity. A variety of mites, beetles, and small moths live as commensals (harmless companions) in the hair, feeding on the algae.

Tree sloths are normally solitary creatures and spend the days and most of the nights sleeping. The little time that remains is spent eating, resting, and moving. This "laziness" contributes greatly to their survival, because their diurnal sleeping posture (hanging from a branch with feet bunched together and head tucked in on the chest) and their coloration give them the appearance of a bunch of dead leaves, rendering them virtually invisible. At night, when they are active, their slow movements do not attract the attention of nocturnal enemies. In addition, sloths have a thick skin and the ability to tolerate strong poisons and severe injuries.

When angered or otherwise disturbed, the three-toed sloth utters a high-pitched cry that sounds like "ai," whereas the two-toed sloth snorts and hisses. Normally, sloths are docile and in the presence of danger remain still, relying on their protective coloration to conceal them. If molested, however, they bite savagely and strike out furiously with the sharp foreclaws.

Burrowing
by
armadillos

Armadillos are primarily nocturnal, spending the daylight hours in their burrows, in expanded nest chambers lined with grasses and leaves. The burrow is also used for the escape from enemies and as a food trap. When using a burrow, in an emergency, some species are able to flex their bony back plates and brace with the feet in such a manner as to become firmly wedged in the tunnel and virtually impossible to pull out. When in danger away from a burrow, armadillos will draw in their feet and lie so that the edges of their armour touch the ground, or will run, burrow rapidly, or even defend themselves with their claws. The three-banded armadillo can roll itself into a ball with feet inward and the head and tail protected by special shields. If captured, an armadillo usually reacts by "playing dead," either stiffening or relaxing, but in either case remaining perfectly still for a brief period. If this course of action does not result in release, the captive armadillo begins kicking vigorously, often facilitating escape. Several kinds of sounds are reported to have been made by fleeing or otherwise agitated armadillos. The peludos (three species of *Chaetophractus*) make snarling sounds, the mulita (*Dasyplus hybridus*) repeatedly utters a guttural monosyllabic sound reminiscent of the rapid fluttering of a human tongue, and other species emit a series of grunts or buzzes when alarmed.

The act of burrowing involves the forelimbs, hindlegs, and tail. The heavy claws on the forefeet are used for scraping and loosening the dirt, pushing it toward the rear. The hindlimbs throw the dirt out behind while the tail elevates the rear end of the animal so that the hindlegs may kick more freely. In this manner, burrows from 0.6 to six metres (two to 20 feet) long, from 0.3 to 1.5 metres (one to five feet) underground and with as many as 12 entrances, may be constructed in a relatively short period. Occasionally, other small animals share these burrows. In North America, rabbits (*Sylvilagus*) and cotton rats (*Sigmodon*) have been found living in the burrows of the nine-banded armadillo.

Armadillos eat insects, mollusks, frogs, small reptiles, carrion, and some plant material and will burrow underneath an animal carcass to obtain maggots and other insects. An acute sense of smell enables some armadillos to detect grubs and other insects up to 12 centimetres (nearly five inches) below the surface of the earth. The peludo is reputed to make conical holes with its snout in the earth in search of food. It first inserts its snout and then rotates its body so that it literally drills the hole. It is able to hold its breath under such circumstances for

up to four minutes. Armadillos have relatively good hearing and vision.

Locomotion. The giant anteater walks with a shuffling gait bearing the weight of its foreparts on the side of the fourth (outermost) digit on the front foot. When hard pressed, this species breaks into a slow, clumsy gallop. It is also a good swimmer. The tamandua walks rather clumsily on the sides of its forefeet when on the ground, in a manner similar to that of the giant anteater. The tamandua is an agile climber, aided both by its foreclaws and hindclaws and by its prehensile tail. Two-toed anteaters are strictly arboreal and are thus very strong climbers. They also have a prehensile tail.

The claws of tree sloths, which are quite similar to those of their terrestrial forebears, have become adapted to a strictly arboreal life. On the ground, a tree sloth is nearly helpless and cannot move at all unless there is something to grasp, then it is able to drag itself along with its claws. Most of the time sloths simply hang upside down, but when they do move they use a slow hand over hand motion both on horizontal and upright branches. Interestingly enough, they are excellent swimmers.

Generally, armadillos walk on the tips of the claws of their front feet and on the soles of the hindfeet; occasionally, they rear up and walk bipedally on their hindlegs. They trot with a stiff-legged gait. At a gallop, some species can outdistance a man. Because of their heavy dermal armour, armadillos have a specific gravity so high that they would normally sink in water. To overcome this problem, they gulp quantities of air that give them the buoyancy needed for swimming. Under some circumstances, armadillos seemingly take advantage of their specific gravity and cross small streams simply by walking under water on the bottom.

Bipedalism
in
armadillos

Habitat. Edentates occur in a variety of habitats, in temperate, subtropical, and tropical environments. Giant anteaters are found in savannas and in wet forests, where they feed either at night or in the day, ripping the terrestrial mounds of termites and ants apart with their powerful foreclaws. Lesser and silky anteaters are both arboreal, the latter totally so. Lesser anteaters prefer tropical forests, parklands, and, to a lesser extent, savannas, in which they forage mostly in trees on ants, termites, and various other insects. They are predominantly nocturnal but may be active at twilight as well. Silky anteaters, which are nocturnal, are found in tall humid tropical forests usually high up in the trees, where they feed almost exclusively on termites. The herbivorous tree sloths occur only in dense tropical rain forests. Armadillos frequent grasslands, parklands, open and even thick forests. They may live largely in subterranean burrows (pink fairy armadillo) in grasslands or dig burrows for temporary shelter in forests (*Dasyplus hybridus*). The shape of the burrow entrance is often diagnostic of the cross-sectional shape of the owner: the nine-banded armadillo (*Dasyplus novemcinctus*), for example, which has a high rounded carapace, digs a nearly circular entrance. The six-banded armadillo (*Euphractus sexcinctus*), which is compressed from top to bottom, digs a burrow with an entrance that is wider than it is high. The habitat preference of the three commonest armadillos in Uruguay illustrates the ecological diversity of this group. *Dasyplus hybridus*, the smallest of the species, prefers grassland and most of its burrows are found there. The six-banded armadillo burrows near streams in open woodlands and grasslands where it occasionally overlaps with the former species and with the nine-banded, which most frequently makes its burrows in woodlands.

Reproduction. The gestation period of those edentates for which it is known ranges from 65 days in the hairy armadillo (*Chaetophractus villosus*) to 263 days in the two-toed sloth. The protracted gestation period of the sloth is attributed to delayed implantation. The very early embryo, consisting of but a few cells, lies free in the uterus without developing further for several months, before becoming imbedded in the uterine wall, at which time normal embryonic growth recommences. Anteaters and tree sloths normally bear single young. Most armadillos give birth to one or two young, although some species demonstrate polyembryony (the birth of from two to 12

Delayed
implanta-
tion

identical young, which develop initially from a single fertilized ovum) and delayed implantation (attachment of the ovum to the uterine wall after a variable delay period). The nine-banded armadillo has both of these phenomena. The delayed implantation is seasonal, the stages occurring at fixed times of the year, related to environmental influences. In the Northern Hemisphere, young of this species are born in March and April. Breeding occurs between June and September, but, regardless of when breeding occurs, implantation is delayed until November and the postimplantation pregnancy is between November and March. Delayed implantation is presumably beneficial because birth does not take place until spring, when food is plentiful.

Little information is available on the birth and early life of many edentates. A recent study of the two-toed sloth (*Choloepus didactylus*) revealed that the young is born as the mother hangs upside down. The baby emerges head first and face upward, and, as soon as its forelimbs are free, it grasps the abdominal hair of the mother and pulls itself up on her chest. The mother may assist in birth by pulling on the young. Subsequently, the female breaks the umbilical cord with her teeth.

The young sloth's eyes and ears are open at birth. It suckles for about five months but does not leave the mother until the ninth month, when adult proportions, but not size, have been attained. Adult size is reached at from two and one half to three years. Young anteaters also cling to their mothers for several months to a year after birth. The giant anteater rides on its mother's back, the two-toed anteater on its mother's tail.

Armadillo young are born with soft leathery skins, which harden as they grow older. In the six-banded armadillo, the eyes do not open for several days after birth, but in the nine-banded armadillo, the young are born with their eyes open. Young armadillos are able to walk when only a few hours old. In the nine-banded armadillo, physical maturity is achieved at six months, but females do not have their first litter until they are two years old. In captivity, the giant anteater has lived 14 years, the two-toed sloth 19 years, the peludo 15½ years, and the three-banded armadillo 11 years.

FORM AND FUNCTION

External appearance. Anteaters (Myrmecophagidae) have elongate heads; tubular muzzles with no teeth; and long, sharp claws. Their toes vary in number from two to four on the forefeet and from four to five on the hindfeet, depending on the species. The thick coat of the giant anteater is gray with lateral white-bordered blackish stripes, and the hair is long and straight, especially on the characteristically stiff bushy tail. Tamanduas, or lesser anteaters, have short hair, which is usually cream-colored to brownish, with a black vestlike area over the upper thorax. The entire underside of the tail, as well as the tip, is naked. The two-toed anteater (*Cyclopes*) has soft silky fur, usually golden yellow in colour. In this species, only the underside of the prehensile tail is naked. Males and females are of approximately equal size in *Cyclopes* and *Tamandua*, but males are much larger in *Myrmecophaga*. All anteaters have enormously developed salivary glands that supply the sticky saliva that covers their long extrusible tongues.

The shaggy pelage of tree sloths grows from the belly toward the back on the body and toward instead of away from the body on the limbs (*i.e.*, downward, when the animal is suspended beneath a branch), thus facilitating the shedding of rain. Tree sloths have pale faces with darker eye rings. The body hair is yellowish to brownish, but the individual hairs are fluted and contain a growth of algae that gives a greenish cast to the entire animal, especially in the rainy season. In the drier conditions of captivity, the algae die and the animal loses the natural coloration. The legs are long, with the forelimbs longer than the hindlimbs, and are designed for suspension of the body rather than for "column-like" support.

Armadillos are characterized by a dorsal bony carapace made up of scapula (shoulder), dorsal, and pelvic shields, separated by a series of movable bands. A bony cephalic

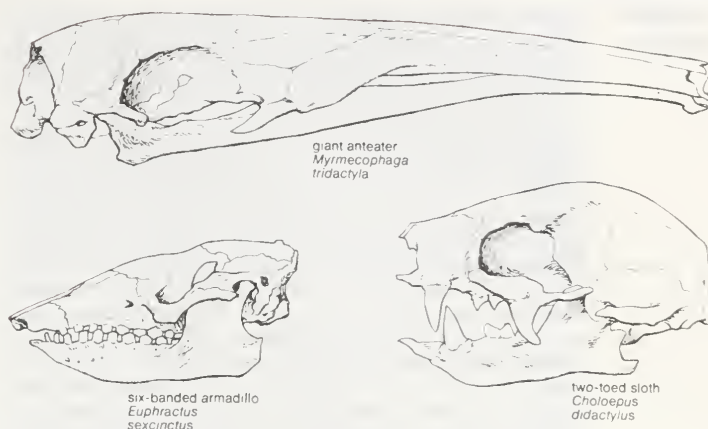


Figure 27: Edentate skulls.

Drawing by J. C. Barbers

(head) shield is also present. The tail is typically armoured as well, except in the naked-tailed armadillos of the genus *Cabassous*. The body hair is greatly reduced, but some, such as the pink fairy armadillo, have dense fur on their sides and underparts. A sparse growth of hair also occurs among the bony plates of all species.

Specializations of internal anatomy. Female edentates have a common urogenital sinus (*i.e.*, both reproductive and excretory functions involve a single chamber) and a simple uterus (a single uterine cavity rather than a branching two-chambered organ). The placenta (the uterine connection between the mother and embryo) is discoid in shape. In other words, the actual areas between which gases, food, and nitrogenous wastes are exchanged are restricted to a single disc-shaped region. The placenta is also deciduous; *i.e.*, it is shed at the time of birth. Females have one or two pairs of mammary glands that are located just lateral to the "armpits," on the chest or on the abdomen. In males, the testes are abdominal and the length of the penis is associated with the degree of difficulty encountered in copulation. In tree sloths, where the partners hang abdomen to abdomen, it is quite small, whereas in armadillos, in which the carapace prevents the partners from getting very close to one another, the penis is long and pendulous.

In general, edentates have poor temperature control. Definitely related to the conservation and manipulation of heat are the extensive anastomoses (fusions) of the veins and arteries (termed *retia mirabilia*) in the limbs and tails of most species. These structures conserve heat, transferring heat from the arterial components carrying warm blood outward along the limbs to the veins lying adjacent for return to the trunk. At the same time, the appendages are cooled.

Tree sloths are heterothermic (*i.e.*, have imperfect control of body temperature) and under experimental conditions their body temperature, normally ranging between 28° and 35° C, may drop to 20° C. At this temperature, the animals become torpid. Temperature control is likewise relatively poor in armadillos. The fact that it builds a subterranean leaf-lined nest has aided the nine-banded armadillo in the northward expansion of its range into the United States. Nonetheless, many individuals of this species perish during severe winters.

EVOLUTION AND PALEONTOLOGY

The ancestry of the order Edentata is believed to be rooted in some unknown but probably generalized insectivore (order Insectivora). The oldest known members of the two presently recognized suborders of edentates were already too specialized for either group to have given rise to the other. Members of the extinct suborder Palaeonodonta, a group restricted to North America, by the time of their first appearance in the fossil record had already developed specializations that eliminate them as possible direct ancestors of the suborder Xenarthra, which flourished throughout the Tertiary in the then isolated South American landmass.

The palaeonodonts, which lived from the late Pale-

The palaeoanodonts

ocene until the beginning of the Oligocene, comprised two families, *Metacheiromyidae* and *Epoicotheriidae*. The most completely preserved form in the former family is *Metacheiromys*, an archaic mammal about 45 centimetres (17.7 inches) in length, reminding one somewhat of an armadillo lacking armour. This animal possessed an assemblage of clearly edentate traits, including the near absence of teeth (except for the enlarged bladelike canines), long compressed foreclaws, and a slight inward twist of the wrist, which is typical of virtually all members of the order. Other metacheiromyids did not show the strong development of the canines. *Palaeoanodon*, known from the upper Paleocene and probably ancestral to *Metacheiromys*, had a series of four or five peglike cheek teeth that lacked enamel. The *Epoicotheriidae*, which had a skull that in outline is reminiscent of that of a small armadillo, were apparently small aberrant forms, highly specialized for a fossorial (digging) livelihood. When the palaeoanodonts passed from the scene at the end of the Oligocene, they left no descendants.

The second suborder, Xenarthra, by far the more successful group of edentates, enjoyed an evolutionary history independent from that of the palaeoanodonts because North and South America were separated throughout most of the Tertiary, not becoming rejoined until the Pliocene. Contemporaneous with the earliest palaeoanodonts of North America was the first South American xenarthran, *Utaetus*, already recognizable as a primitive armadillo and thus sufficiently advanced to eliminate any possibility that known palaeoanodonts might be directly ancestral to xenarthrans. The fact that primitive armadillos are the first known members of the suborder strongly suggests that the other xenarthran families are in some way derived from the Dasypodidae.

The Tertiary radiation of edentates in South America resulted in several more or less separate lineages, which are conveniently arranged in three infraorders. The first of these infraorders, Loricata, includes the armadillos (*Dasypodidae*) and their relatives, the extinct glyptodonts (*Glyptodontidae*) or so-called turtle armadillos, distinguished from the former especially by their possession of a rigid carapace. Armadillos flourished in South America in the Tertiary and were abundant by the Miocene. In the Pleistocene, some species had moved northward as far as what is today the central United States.

Fossil armadillos

Four distinctive lines developed in the Dasypodidae, comprising the subfamilies Dasypodinae, the Stegotheriinae, the Pampatheriinae, and the Chlamyphorinae (including only the living genera *Chlamyphorus* and *Burmeisteria*). The Dasypodinae includes most of the living genera of armadillos as well as *Utaetus*. The fossil members of this group varied in size from that of a small contemporary armadillo to animals as large as the South American Pliocene form *Macroeuphractus* that had a skull nearly 28 centimetres (11 inches) long. The head and body of the latter must have been almost two metres (6.6 feet) long. The fossil members of this subgroup all seemed to favour invertebrates in their diets, but they probably ate carrion and some plant material as well. This distinguished them from the extinct Stegotheriinae, a subfamily whose members were highly specialized for a strictly insectivorous diet. *Stegotherium* of the lower Miocene of South America had an elongate anteater-like skull with much reduced cylindrical cheek teeth.

The Pampatheriinae, known from the late Miocene to the Pleistocene, was a group characterized by great size. At least one form, *Pampatherium*, a nearly complete skeleton of which was recovered from a Pleistocene deposit in Texas, was nearly as large as a rhinoceros. This animal and its closest relatives were apparently adapted to an herbivorous diet. Both insect- and plant-eating groups are thus included among the true armadillos.

The genus *Peltephilus* and several other closely related genera of aberrant armadillo-like creatures have been placed in their own family, Peltephilidae, because, among other characteristics, they possessed at least one pair of hornlike protuberances on the muzzle. These animals, which appeared in the Oligocene and disappeared in the Pliocene, were apparently specialized for flesh eating.



(Left) Ground sloth (*Megatherium*) and (right) glyptodonts (*Doedicurus*). Painting from fossil skeletons by Charles R. Knight.

By courtesy of Field Museum of Natural History, Chicago

The glyptodonts (*Glyptodontidae*) evolved from primitive armadillos, first appearing in the fossil record in the Eocene with the genera *Glyptatelus* and *Lomaphorelus* but remaining relatively uncommon until the beginning of the Miocene. Subsequently, they experienced a radiation of their own and by the Pleistocene several kinds had spread into North America. Throughout their history, they demonstrated an increase in size culminating in the Pleistocene giant, *Doedicurus*. They became extinct before the end of that epoch.

The giant glyptodont

The second infraorder, Pilosa (hairy) includes the sloths, both the living arboreal forms (*Bradypodidae*) and the fossil ground dwellers of the families *Megatheriidae*, *Megalonychidae*, *Myodontidae*, and *Orophodontidae*. There are no fossils of tree sloths, but they appear to be structurally related to the ground sloth families *Megatheriidae* and *Megalonychidae*. The oldest fossils of ground sloths belong to the *Megatheriidae* and come from Oligocene deposits in South America. By the early Miocene, primitive members of this group were rather common. One form, *Hapalops*, although typical of that time, was of relatively small size for a ground sloth. It had a slender body about 1.3 metres (4.5 feet) in length and had five toes with developed claws on both its forefeet and hindfeet. *Hapalops* had a moderately elongated skull with slender premaxillae (bones of the upper jaw), and a process resembling a spout that protruded from the lower jaw and may have supported horny plates for cropping. This animal had five upper teeth and four lower teeth in each jaw, the first upper pair of which resembled canines. The small Pleistocene ground sloth *Nothrotherium* apparently was a descendent of *Hapalops* and the two comprise a branch of the *Megatheriidae*, which throughout its history was given to small size. This stock differed greatly from the main stem of the family, which had no specialized "canines" and had the two inner toes reduced, represented by the Pleistocene *Megatherium*, a massive ground sloth larger than an elephant.

The megalonychids, although related to the megatheriids, were distinguished from them by their definite canine-like teeth and lack of the spoutlike process on the lower jaw. The megalonychids ranged in size from one island-dwelling animal no bigger than a domestic cat to the mainland ox-sized *Megalonyx*. Both the *Megalonychidae* and the *Megatheriidae* reached southern North America by the Pleistocene, but they became extinct on the Continent at the end of the epoch. Some megalonychids invaded the West Indies, however, and experienced a minor radiation in Hispaniola, Puerto Rico, and Cuba. One or more species may have flourished there to within a few hundred years of the present.

Distinct from the megatheriid-megalonychid line of ground sloths were the *Myodontidae*, which first appeared in the Oligocene; from the Pliocene until their extinction at the end of the Pleistocene, they were the dominant South American group of ground sloths. The myodonts

Mylodonts differed from the other ground sloths in their partially developed upper canines, reduced first toe of the hind-foot and rounded ossicles (bony tubercles) in their skin. The mylodonts followed the trend of other ground sloth families toward great size and awkwardness, which in this group culminated in the Pleistocene giant, *Glossotherium*, which was the size of a large ox. As with other xenarthran groups in the Pleistocene, the mylodonts extended as far north as the southwestern United States.

The other family of ground sloths, the Orophodontidae, is known only from the Oligocene of South America. This briefly-lived side group is represented by *Orophodon* and *Octodontotherium*, both of which were rather small sloths that do not seem to have left any descendants.

The final infraorder, Vermilingua, comprising the anteaters (Myrmecophagidae), has not been found in the fossil record any earlier than the Miocene. The fossil myrmecophagids, as represented by the Pliocene form *Palaeomyrmedon*, were already completely toothless and presumably well adapted to an insectivorous diet. The Myrmecophagidae first appeared in South America in the early Miocene and apparently only fairly recently have spread to Central America.

CLASSIFICATION

Distinguishing taxonomic features. The degree of development of certain sets of structural characteristics taken in concert defines the Edentata and its subgroups. None of these features is universally distributed throughout the order. The most important are as follows: (1) in addition to the normal zygapophyses (protuberances), the occurrence of supplementary articular surfaces called xenarthrous articulations between the arches of sequential posterior trunk vertebrae; (2) development of certain projections of the scapula bone, or shoulder blade, called the acromion and coracoid processes; (3) articulation of the ischium bone of the pelvis with the proximal caudal (tail) vertebrae, forming an elongate sacrum; (4) length and massiveness of limb bones; (5) fusion of distal limb bones; (6) hyperdevelopment of claws; (7) development of forefeet and hindfeet; (8) a skull characterized by a small brain housed in a long tubelike brain case, typically with a reduced premaxilla (*i.e.*, elongate palate), often an incomplete zygomatic arch, and a reduction in number of teeth and their enamel covering.

Annotated classification. The classification presented here combines those proposed by the two paleontologists R. Hoffstetter and Alfred S. Romer. The features given below are only those that distinguish each group from the others. Groups indicated by a dagger (†) are known only as fossils.

ORDER EDENTATA

Placental mammals with reduced dentition, often lacking enamel; heavily clawed forelimbs for burrowing. Coracoid process of the scapula more strongly developed than in other eutherian mammals. About 31 recent species.

†Suborder Palaeanodonta

Fossil only. Bony carapace and the xenarthrous articulations on the vertebrae lacking; no fusion between the ischium and caudal vertebrae; simple teeth, usually single-rooted and lacking enamel; slight inward twist to the forelimbs.

†*Family Metacheiromyidae.* Upper Paleocene to lower Eocene. Small; lightly built edentates with sharply pointed somewhat laterally compressed canines; postcanine teeth reduced or lacking; armadillo-like feet with 4 to 5 digits; the lateral digits reduced, the inner greatly so. North America.

†*Family Epopoitheriidae.* Lower Eocene to lower Oligocene of North America. Very small burrowing edentates, with a domed skull having a depressed snout and possessing a zygomatic arch; having 6 slightly differentiated cheek teeth, nearly cylindrical, lacking enamel, and single-rooted.

Suborder Xenarthra

Some extinct and all living edentates; all with xenarthrous posterior trunk vertebrae. Lack teeth or have only poorly differentiated cheek teeth, lacking enamel and having single roots. Fusion (symphysis) between the ischium and the anterior caudal vertebrae in all except the 2-toed anteater (*Cyclopes*).

Infraorder Loricata

Bony plates covered with horny plates, forming a protective carapace over the trunk.

Family Dasypodidae (armadillos). Lower Eocene to Recent in South America; Pleistocene to Recent in North America; forests, savannas, and plains; primarily in the tropics and subtropics. Teeth more than $\frac{5}{6}$ (6 above and below on each side); top and sides of the body covered with horny scutes over a bony, flexible carapace.

†*Family Peltephilidae.* Upper Oligocene to lower Pliocene (questionably) in South America. Teeth $\frac{7}{2}$; short, broad skull; hornlike structures on the head armour.

†*Family Glyptodontidae* (glyptodonts). Mid-Eocene to Pleistocene in South America; upper Pliocene to Pleistocene in North America. Fusion of bony plates on the back into a solid turtle-like carapace; teeth $\frac{5}{2}$ with 3-lobed pattern.

Infraorder Pilosa

Incomplete zygomatic arch (bony ridge below eye); 4 or 5 simple cheek teeth; expanded acromion, which unites with the coracoid region.

†*Family Orophodontidae.* Oligocene in South America. Bilophodont teeth, $\frac{5}{4}$, made of a mass of compacted dentine; highly specialized astragalus (ankle bone). Small to medium body size.

†*Family Megalonychidae* (ground sloths). Upper Pliocene to Pleistocene in South America; middle Pliocene to Pleistocene in North America; Pleistocene in West Indies. Ground sloths. Large canine-like teeth in the upper and lower jaws and no spoutlike terminus on the lower jaw.

†*Family Megatheriidae* (ground sloths). Upper Oligocene to Pleistocene in South America; Pleistocene in North America. Partial development of upper "canines"; spoutlike terminus on the lower jaw.

†*Family Mylodontidae* (ground sloths). Upper Oligocene to Pleistocene in South America; Pleistocene in North America. Distinguished by a partial development of upper "canines," triangular-shaped cheek teeth, and a hindfoot with the 1st toe reduced.

Family Bradypodidae (tree sloths). South and Central America; tropical forests; no fossil record. Dental formula of $\frac{5}{4,5}$; 6 to 9 cervical vertebrae.

Infraorder Vermilingua

Toothless and hairy edentates.

Family Myrmecophagidae (anteaters). Lower Miocene to Recent in South America; Recent in Central America; tropical forests and savannas. Teeth lacking; tubular mouth with a small terminal opening.

Critical appraisal. The edentates are a natural assemblage of several seemingly rather diverse groups. Given that the oldest known representatives of the palaeonodons and the xenarthrans are too advanced for either to be ancestral to the other, it would be of great interest to find the common ancestor of both groups, as well as the insectivore precursor of this ancestor. To date, no good fossil candidates have been unearthed in either category.

Authorities disagree on the relationships and degree of distinctness of some groups of fossil edentates. The palaeonodons are believed by some to be closer to the Old World pangolins (order Pholidota) than to the New World edentates. The armadillo-like *Pseudophorodon* from the upper Oligocene of South America is sometimes placed in its own family but is considered here to belong in the Peltiphilidae. The palaeopeltids, known especially from polygonal plaques of bone, are considered by some to have been aberrant glyptodonts but are considered by Romer to have been the distinct superfamily Palaeopeltoidea. The unusual Miocene ground sloth *Entelops*, due to the presence of teeth on the premaxillary bone, is sometimes separated from the Mylodontidae as sole member of the family Entelopsidae.

The degree of relatedness of sloths and anteaters has been the subject of disagreement. It has been shown that despite certain adaptations of the front limbs and digestive tract for insect eating, anteaters are close to the most primitive xenarthran edentates in many characteristics of the forefeet and hindfeet. Hofstetter has suggested that the evolutionary lines that produced the modern Myrmecophagidae and Bradypodidae separated early and that, therefore, the two families are not closely related but should be placed in separate suborders.

The origin of the tree sloths and the relationship of the currently recognized two genera is a question of some

Relation-
ship of
sloths and
anteaters

interest. As mentioned above, the two genera may be independently derived from two extinct families of ground sloths. *Choloepus* shows affinities to the megalonychids, and *Bradypus* is more like the megatheriids.

At one time, the Old World scaly anteaters or pangolins (order Pholidota) and the South African aardvarks (order Tubulidentata) were included in the Edentata. The adaptations for an insectivorous diet that these two groups seem to have in common with the New World anteaters and the armadillos are most assuredly a result of convergent evolution and not a reflection of any close phylogenetic relationship. (J.C.B.)

Lagomorpha (rabbits, hares, pikas)

The terrestrial mammals of the order Lagomorpha include the relatively well-known rabbits and hares and also the less frequently encountered pikas, or mouse-hares. Rabbits and hares characteristically have long ears, a short tail, and strong hindlimbs that provide a bounding locomotion. In contrast, pikas have shorter, rounded ears, no external tail, and less well-developed hindlimbs associated with scampering locomotion.

GENERAL FEATURES

Wild lagomorphs are small to small-medium in size, ranging from the smallest pikas, about 150 millimetres (5.9 inches) in length and 100 grams (3.5 ounces) in weight, to the largest hares, 700 millimetres (27.6 inches) and 4.5 kilograms (10 pounds). Wild rabbits range between pikas and hares in size, while some varieties of domestic rabbit may reach up to seven kilograms in weight.

Nearly worldwide in distribution, lagomorphs are absent only from most of the Southeast Asian islands, Australia, New Zealand, Madagascar, southern South America, and Antarctica. Humans have introduced rabbits and hares into areas outside their original ranges for purposes of sport and to provide a readily available food supply. The natural range of the Old World or European rabbit, *Oryctolagus cuniculus*, appears to have been southwestern Europe and North Africa; however, in Roman times this species was often introduced to islands, where it flourished, and spread north and east, partly following human agricultural activities, reaching Britain in Norman times. More recently, *Oryctolagus* was introduced for sport and food into New Zealand and Australia, where the absence of native predators and other factors soon led to its becoming an agricultural pest and a threat to some of the native fauna. In North America several species of the cottontail rabbit, *Sylvilagus*, have shown an increase in abundance and a broadening of range in areas disturbed by human activity and settlement.

Wild lagomorphs are popular with hunters for sport as well as for food and fur. Domestic rabbits, all descendants of *O. cuniculus*, are raised for meat and skins, the latter of which are used both as pelts and for making felt. Domestic rabbits make good and relatively undemanding pets. Their attractive appearance and quiet manners have made them favourite characters with children's storytellers. On the other hand, wild rabbits and hares locally may become pests, depleting vegetation available to domestic grazers or damaging young trees and orchards, especially in winter. The problem becomes especially acute in areas where human populations have removed the natural predators of lagomorphs. Attempts have been made to control Old World rabbit populations in western Europe and Australia by introducing the viral disease myxomatosis, which exists naturally in populations of certain South American rabbits of the genus *Sylvilagus*, but in some areas the rabbits have developed resistance to the disease. Use of this virus in controlling rabbit numbers has been fought by some humane organizations because of the distressing symptoms of the disease. Pikas, usually found in areas remote from human activity, have little economic importance.

NATURAL HISTORY

Rabbits and hares occupy a wide variety of habitats, including grassland, desert, forest, marsh, brushland, and tundra. They are found in mountains up to elevations of

4,900 metres above sea level. Rabbits frequent areas where cover is readily accessible and do not venture far from it. Some rabbits dig their own burrows, whereas others use burrows dug by other mammals. Although most rabbits are solitary, some live in pairs. The Old World rabbit, *O. cuniculus*, usually is gregarious, forming breeding colonies known as warrens, with extensive burrow systems. The Mexican volcano rabbit, *Romerolagus*, maintains runways through dense vegetation. Although rabbits are typically terrestrial, two North American species of *Sylvilagus*, the marsh rabbit and the swamp rabbit, take readily to water when necessary. Hares tend to live in more open areas than do rabbits. They are solitary and do not use burrows but rest and take shelter in shallow depressions made in soil or vegetation. Pikas are found in northern steppes, semideserts, some forests, and scrub thickets of Asia and in rocky terrain, especially on rock debris on slopes; they range in mountainous areas to 6,000 metres (19,700 feet) in northern and central Asia and western North America. Steppe-dwelling pikas construct and inhabit burrows in flat areas. Rock-dwelling pikas usually seek shelter among rocks and boulders. Some pikas are colonial, whereas others are solitary and may exhibit some degree of territorial behaviour.

Rabbits and hares are customarily nonvocal, making cries or screams only when frightened or injured. One rabbit, the South African red "hare" *Pronolagus*, is known to utter a warning call. Pikas are highly vocal, having a whistle or bark and a chattering call, reflected in such common names as whistling hare and piping hare. The voice is used in giving alarm signals and in maintaining territorial boundaries. Mountain hikers or climbers frequently hear their calls and may see these attractive little creatures sunning themselves on rocks and going about their daily activities.

Rabbits and hares are active mainly from dusk to after dawn, whereas most pikas are active during the day. All lagomorphs are plant eaters, grasses and soft vegetation being their dietary staples, but they may also eat twigs and bark, especially in winter. Arctic hares are known to break through crusted snow with their forelimbs to feed on the buried vegetation. Rabbits and hares are not known to store food, but pikas store and cure piles of vegetation for use during the winter. Some pikas cut grasses, carry them to the pile, and allow the vegetation to dry before adding more, whereas other pikas actually turn the contents of a pile to aid in the curing. Pikas that live among rocks first pile the vegetation in the sun and open air before moving it under the shelter of rocks for the winter. Steppe-dwelling pikas pile their hay in the open. By contrast, however, other pikas may store hay not in piles in the open but in crevices in trees, stumps, and among rocks. Some of the Asian pikas unwittingly supply hay for domestic livestock during the winter.

Some lagomorphs regularly reingest fecal pellets (coprophagy). Two kinds of pellets are produced: dry and hard, and soft and moist. The latter variety, which appear to contain both vitamins and metabolic products, are eaten, in most cases directly from the anus. The nutritional effect of this practice has been compared to that of rumination among cows.

Rabbits and hares differ from one another in the condition of their young at birth. Young rabbits are born and cared for in a nest; they are naked, blind, and relatively helpless at birth. Hares are born in the open; they are furred, have open eyes, and can run shortly after birth. Confusion of the terms "rabbit" and "hare" in common parlance is shown by such popular names as "jack rabbit" for a hare and "Belgian hare" for a variety of domestic rabbit.

The reproductive abilities of rabbits are proverbial. Rabbits and hares usually produce several litters during each breeding season, two or three litters being common among hares and three to six among rabbits. The gestation period ranges from about 28 days in rabbits to 47 in hares, and litter size is usually between two and eight. A female cottontail rabbit from an early spring litter can produce offspring by the end of the summer at age 10 weeks. In contrast to most mammals, female rabbits and hares are

Vocalizations

Economic importance of lagomorphs

Reproductive potential

usually larger than males. Reproductive cycles of pikas are less well known than those of rabbits and hares. Their breeding season occurs in late spring and summer, with a gestation period of approximately 30 days. Litter size normally ranges from two to six, and two or three litters may be born each year. In some pikas, the newborn young are only lightly furred, have eyes and ears closed, but can move about after about eight days.

Lagomorphs, especially hares and rabbits, are important elements in many food chains. They are preyed upon by a variety of carnivores, both mammals and birds, who rely upon them as dietary staples. Wolves, foxes, bobcats, weasels, predatory hawks, and owls all take their toll. More northerly species of rabbits and hares frequently exhibit cyclic fluctuations in numbers, which are reflected in the cycles of the carnivores dependent upon them. The snowshoe hare-snowy owl cycle is one of the prominent ones. Pikas are preyed upon by various carnivores, among which weasels and birds of prey are probably most important.

Externally most rabbits and hares are not dissimilar from one another, the main variation occurring in length of ear and strength of hindlimb. The coat is usually brownish or reddish brown above and lighter to white below. An exception to this coloration is found in the Arctic hare, in which the winter coat is pure white with only the ear tips dark. In more northerly representatives of the Arctic hare, individuals retain this colour pattern all year, whereas more southerly forms of it may become grayish or buff-coloured. In other northern hares, the varying hare (*Lepus timidus*) of Eurasia and the snowshoe hare (*L. americanus*) of North America, the winter coat is white and the summer coat brown. Some varying hares in Ireland and snowshoe hares from the Pacific Coast do not become white in winter. The most unusual coloration in a modern rabbit occurs in the Sumatran rabbit *Nesolagus*, which has brown stripes on its buff-gray pelage and a bright red rump and tail. Pikas are usually brownish to reddish above and somewhat lighter below.

Lagomorphs illustrate adaptation to particular habitats in an interesting gradational series of features for detection of, and escape from, enemies. At one end of the series are the long-eared, wary, and alert hares, inhabitants of open country, which detect enemies at considerable distance and rely for escape on their strengthened hindlimbs and bounding locomotion that takes them up to 80 kilome-

tres (50 miles) per hour. Rabbits, shorter eared and with weaker hindlimbs, do not react to such distant threats and are scamperers or bounders who do not venture far from cover. They are quick over short distances but lack the endurance of hares. With still shorter ears and relatively weaker hindlimbs, pikas live in areas in which a quick scamper takes them the short distance that usually separates them from secure cover.

FORM AND FUNCTION

Lagomorphs are well adapted to a herbivorous diet. Like rodents, they have well-developed incisors that grow continuously from the roots, while they are worn down at the cutting edges. These teeth are extremely effective for severing plant stems and for gnawing on bark. The very mobile lips each consist of two lobes, which meet behind the chisel-like incisors when the mouth is closed. There are no canine teeth, a gap (diastema) showing where they might have been. Cheek teeth farther back in the jaw are also ever growing, wearing away as they grind abrasive vegetation. The upper tooth rows are more widely separated than the lower rows, and chewing is done with a transverse movement. The chewing muscles of the jaws, though strong, are less well developed than in rodents.

Vegetation passes through the long small intestine, which has a spiral valve (like a wood screw), providing a much greater surface area for digestion of food and absorption of nutriment. A large pouch, or cecum, located at the point of attachment of the large intestine contains bacteria that aid digestion and that produce the nutritional soft fecal pellets.

Skeletal adaptations for speed and agility are evident and are discussed as taxonomic characters in the *Annotated classifications* section, below. The bones of the hindlimb are fused where they move against the anklebone (calcaneum), affording strength and leverage employed in bounding or scampering locomotion. Lagomorphs move about as if they walked on their toes (digitigrade locomotion). There are five digits on the forefoot and four or five on the hindfoot.

EVOLUTION, PALEONTOLOGY, AND CLASSIFICATION

The order Lagomorpha appears to have originated in northern Asia, probably by the end of the Paleocene (about 55,000,000 years ago), and has been relatively

Dentition

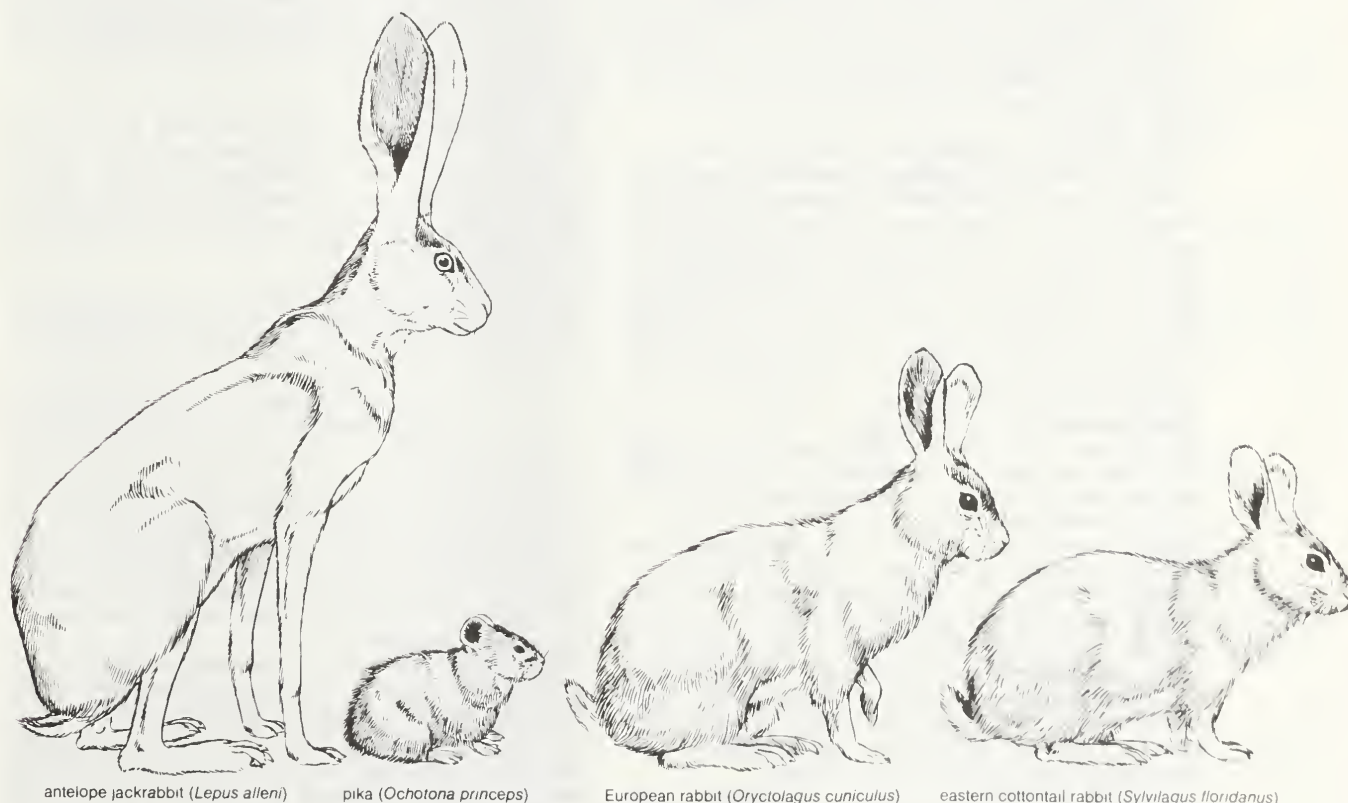


Figure 28: Diversity among lagomorphs.



The summer fur (left) of the snowshoe hare (*Lepus americanus*) blends with the browns of woodland and prairie. A mottled effect, created by white hairs gradually replacing the brown, is similar to patches of new snow, while the white winter coat (right) is completely camouflaged in a snowy landscape.

The National Audubon Society Collection/Photo Researchers (left) Charles J. Ott (right) Leonard Lee Rue

Evolutionary trends

stable morphologically throughout the approximately 40,-000,000 years since the end of the Eocene, when its fossil record first becomes well documented. Lagomorphs seem to have established early in their history the patterns of dental and skeletal development that were to change, for the most part, only by gradual advances in the basic pattern. Within the order several trends can be observed from more primitive to more advanced lagomorphs: development of ever-growing cheek teeth, simplification of the pattern of the occlusal surfaces of the upper cheek teeth, and shortening of the shaft of the lower incisor where it passes back into the jaw. When dealing with fossil forms, often known only from fragments of jaws and teeth, the most useful characters that differentiate taxa are ordinarily found in the structure of the anterior two upper, and anterior lower, premolars. Presence or absence of the last upper molar may also be important, as may the pattern of the upper cheek teeth. When preserved in fossils, structure of palate and of lower jaw also provide characters that assist in determining affinities.

The family of rabbits and hares (Leporidae) entered North America by the end of the Eocene and underwent most of its Middle Tertiary development there. By the Pliocene (about 7,000,000 years ago) it had become reestablished in Asia and had moved also into Europe. The leporids now extend throughout those ranges and down to the tip of South Africa and as far south as northern Argentina in South America.

The pika family (Ochotonidae) spread from Asia to Europe, where they developed into a number of types in the Middle Tertiary. One of these lines persisted until the late Pleistocene or early modern times (about 1,000,000 years ago) on Corsica and Sardinia. Other ochotonid branches reached Africa and North America in the Middle Tertiary. The extant genus *Ochotona* appeared in Asia in the Pliocene and spread from there, reaching western Europe and eastern North America in the Pleistocene. The current range of the genus represents a considerable reduction from that during the Pleistocene.

Distinguishing taxonomic features. Within the order Lagomorpha living members of its two extant families are clearly separable on morphological grounds, the differences being in skull, jaw, and dental structures, as well as in postcranial morphology.

Morphological trends in the rabbits and hares include increase in arching of the skull, correlated with the development of bounding locomotion and of a relatively upright posture of the head; reduction in the ability to perform lateral movements, caused partly by fusion of bones and partly by modifications related to specialization for anteroposterior movements; strengthening of the hindlimbs and pelvic girdle; and elongation of the limbs distally. Pikas lack these skeletal modifications and are adapted rather for a scampering locomotion.

Annotated classification. In older classifications lagomorphs were usually placed in the order Rodentia. Both

rodents and lagomorphs are gliriform mammals; that is, they have enlarged, ever-growing incisor teeth that function as effective tools for gnawing. Under this scheme, the lagomorphs were the suborder Duplicidentata, and the true rodents the suborder Simplicidentata. Early workers, however, recognized many features separating these suborders. For one, there are two pairs of upper incisors in lagomorphs (thus the name Duplicidentata) as opposed to one pair of upper incisors in rodents (Simplicidentata). It is now known that lagomorphs and rodents have long separate histories and that even their earliest representatives did not closely resemble one another. Accordingly, it is more accurate and currently generally accepted to classify the two groups as distinct orders, Lagomorpha and Rodentia. The following classification is the one in general use.

ORDER LAGOMORPHA

Small to medium-sized mammals with two pairs of upper incisor teeth, one pair large, ever-growing, and grooved, in front of a pair of small, peglike incisors, and one pair evergrowing lower incisors. In skull the jawbone, or maxilla, pierced or with open lacework laterally; palatal bone with large, deep-cut openings and relatively short bridge. Of muscles used in chewing, masseter muscles not greatly enlarged, temporalis muscle reduced, pterygoid muscles well developed. Twelve extant and approximately 35 known extinct genera are included in three families, one extinct (†).

†Family Eurymyidae

Late Paleocene and (?) Eocene. This extinct family of two genera known only from northern Asia may or may not be closely allied to the other two families of lagomorph. Several anatomical features, such as the presence of only two pairs of upper premolars, show that it was not directly ancestral to later lagomorphs, but it may have been near a group from which later lagomorphs originated.

Family Leporidae (rabbits and hares)

Late Eocene to Recent. Almost worldwide. Distinguished by elongated ears, hindlegs much longer and stronger than forelimbs, often greatly strengthened, and short external tail. Most anterior premolars unlike others in pattern, the following two premolars similar; usually three upper and three lower molars in each jaw, reduced to two upper molars in *Pentalagus*; shafts of upper cheek teeth extend upward into socket in bony capsule (orbit) surrounding eye. Skull with jawbone having open lacework of bone laterally; bony process present over orbit. The palate has one pair of large deep-cut holes; bony palate with longer maxillary and shorter palatine components. Tendency toward arching of skull in profile, related to increasingly upright posture of head. Clavicle rudimentary. Strong, well-developed lumbar region of back, and strong pelvic girdle for articulation with hindlimb and attachment of limb muscles. Five digits on forefoot and on hindfoot.

The subfamily Palaeolaginae is an extinct group that includes primitive Old and New World forms. Descendants of the Palaeolaginae were the extinct Archaeolaginae, including most North American Middle Tertiary forms as well as a group that reached Eurasia in the Pliocene. The most advanced leporids are in the subfamily Leporinae, which includes all living leporids (11 genera), including the American *Sylvilagus*, the Old World *Oryctolagus*, *Lepus* of the Northern Hemisphere,

Romerolagus of central Mexico, *Pentalagus* of the Ryukyu Islands south of Japan, *Bunolagus* and *Pronolagus* of South Africa. Other leporids showing some relatively primitive characters are *Nesolagus* of Sumatra and *Brachylagus* of the western United States. The subfamilial divisions are based partly on the pattern of the most anterior upper and lower premolars, which seem to be valid indicators of leporid relationships. Size 25 to 70 centimetres (9.8–27.6 inches).

Family Ochotonidae (pikas)

Oligocene to Recent. Asia and western North America. Distinguished by short, rounded ears, hindlegs slightly longer than forelimbs, no external tail. Three upper premolars exhibit three different patterns; two upper and three lower molars in each jaw in modern forms, ranges from three upper and lower molars in each jaw to two in each jaw in fossil forms. Shafts of upper teeth curve outward into arch of the zygomatic bone beneath eye socket. Skull with single large opening in each jawbone, other holes present in some fossil forms. Bony process over orbit absent in modern ochotonids but present in some fossils. On palate one or two pairs of deep-cut holes; bony palate with larger palatine component and smaller maxillary component. Skull typically more flattened than in leporids. Well-developed clavicle. Five digits on forefoot and four on hindfoot. One extant genus, *Ochotona*, and 14 or 15 extinct genera. Size 15 to 28.5 centimetres (5.9–11.2 inches).

Critical appraisal. Perhaps the most important unanswered question in lagomorph classification is that of affinities to other mammals. Part of the difficulty is that the fossil record from the Early Tertiary of Asia is incomplete for large parts of the Paleocene and Eocene. The Eurymylidae, known from some Paleocene and possibly Eocene deposits in Asia, have some characters similar to those later lagomorphs but other characters showing that they were not the direct ancestors of later forms. Other fossil mammals now known from the Paleocene of Asia may represent forms related to eurymylids. The evidence is incomplete, but perhaps eurymylids and these possibly allied forms were in a generally ancestral position to leporids and ochotonids. The wider affinities of this possibly ancestral group to other mammals are not yet clear. No transitional forms are known between the eurymylids and later lagomorphs, which had attained their distinctive characters by the late Eocene. Earlier attempts to arrive at the solution to the problem of lagomorph affinities often emphasized similarities between various primitive herbivorous animals, such as the condylarths, and primitive lagomorphs. These similarities, however, are probably due to the presence, in both groups, of morphological characters common to primitive mammals in general.

A further point causing some difficulty in studies on fossil leporids and ochotonids is the assignment of various primitive forms to one family or the other. The problem is especially acute when the fossil record is based only on a few teeth or jaw fragments. Morphologically these primitive forms seem to be intermediate between the two families.

Some differences of opinion exist in regard to the subfamilial classification of leporids. Some students of these animals have taken the step of subdividing leporids into several families. Other problems in the classification of modern leporids are matters of specific assignment of various populations. A sound subfamilial arrangement for the ochotonids has not yet been established. (M.R.D.)

Rodentia (rats, mice, beavers, squirrels, guinea pigs, capybaras)

The rodents, or Rodentia, are the most abundant order of mammals. At present, over a quarter of the families, 35 percent of the genera, and 50 percent of the species of living mammals are rodents. Probably an even higher percentage of individuals are rodents, for they tend to be small animals with dense populations. They are one of the few groups of animals that flourish in close association with men. Some, such as squirrels, live independently but fairly successfully near humans. Others, such as the house mouse (*Mus musculus*) and black and Norway rats (*Rattus rattus* and *R. norvegicus*), have adapted themselves to human civilization, and live everywhere that man does. These two rats (and the Polynesian rat, *Rattus exulans*, of Australia and Oceania) have travelled in ships and boats of all sizes, and have pop-

ulated the entire habitable world, especially near human habitations.

GENERAL FEATURES

All rodents possess one pair of upper and one of lower incisors, growing throughout life, with the enamel restricted to a band on the front side of the teeth. Behind this is a large gap (diastema) followed by two to five cheek teeth. The jaw articulation is so arranged that when the cheek teeth are in use, the incisors do not meet, and vice versa. The incisors grow continuously, and must be worn off equally fast, or the whole gnawing mechanism is ruined. Because of the necessity to abrade these incisors, rodents spend a considerable amount of time gnawing hard objects.

Generally rodents are small. Some mice and dormice are among the smallest of living mammals, adults being as small as 75 millimetres (three inches) long, including the tail, and weighing as little as 20 grams (0.7 ounce). The largest living rodent is the South American capybara (*Hydrochoerus hydrochoeris*), reaching over 1.3 metres (four feet) in length and as much as 50 kilograms (about 110 pounds) in weight. A fossil rodent recently described from Uruguay is reported to have had a skull as large as that of a bull and a body bulk as large as that of a wild boar.

Rodents are of major economic importance, primarily as consumers of the grains that are the basic foodstuff for

Size range

Relation-
ships with
other
mammals



Figure 29: Range of body-plan variation of Rodentia.

man. It has been estimated that rats and mice destroy up to one-third of grain crops under conditions of heavy infestation. Burrowing rodents may damage root crops. The muskrat (*Ondatra zibethica*) and nutria (*Myocastor coypus*), introduced into Europe as fur sources, have escaped and spread over much of Europe between the Baltic and the Alps. Their burrows, particularly in canal banks, have been a major source of damage to the drainage system, most especially in The Netherlands. A number of rodents serve as reservoirs for human diseases, such as bubonic plague, tularemia, scrub typhus, and others. The plague that ravaged Europe during the mid-14th century was transmitted by fleas from rats to humans.

Several rodents (beaver, muskrat, chinchilla, nutria, squirrel) produce fur useful to man. All but beaver and squirrel have been domesticated for this purpose. Albino mice and rats, hamsters, and guinea pigs are widely used as laboratory animals for biological and medical purposes. Guinea pigs were domesticated by the Incas for food; a few kinds of rodent have been raised as pets.

Distribution of rodents

Rodents occur naturally in all parts of the land where there is an adequate food supply and are found in essentially all terrestrial habitats. They range from well above the Arctic Circle to the southern tips of Africa and South America and were the only terrestrial placental mammals, other than bats, to reach Australia before the arrival of man. Many rodents have successfully adapted to difficult environments such as deserts. Many rodents have broad climatic tolerances, an example being the North American porcupine, which is found from the Arctic Circle to central Mexico and from the Atlantic to the Pacific. Most are quadrupedal scampers, but they generally have much freedom of use of their forefeet in manipulating food; many are burrowers, spending most of their life underground; some are ricochetal, leaping on their hind legs; flying squirrels use skin membranes to glide from one tree to another; a few (beaver, muskrat, water vole, nutria) have become amphibious in habit, living in freshwater streams

and ponds; and a number of South American rodents are cursorial (running) animals.

IMPORTANCE TO MAN

Destruction of crops and foodstuffs. The most important rodents, from the point of view of economic damage, are the Norway and black rats, with the house mouse close behind them. It has been frequently estimated that the rat population of the United States is approximately equal to the human population. A population of over 1,000 rats per acre (2,470 per hectare) on an Iowa farm has been reported. Population explosions of house mice occurred in the Central Valley of California in 1926-27 and 1941-42. During the former, the mouse population was estimated to have reached over 80,000 per acre (198,000 per hectare).

Population explosions

Remarkable population explosions of voles (*Microtus*) and wood mice (*Apodemus*) are well known in western Europe, recurring every few years. In France from 1790 to 1935 there were at least 20 mouse plagues, some lasting several years. Estimates of abundance are highly inaccurate, but there are reports of 8,000 voles per acre (19,800 per hectare) and 15 to 20 vole burrows per square metre (13 to 17 per square yard) in peak conditions. There was a population explosion of voles and wood mice over much of Germany beginning in late summer of 1917 and lasting through 1918. Damage to crops was serious: clover was the favourite food; winter rye and wheat, sugar beets, and potatoes also were destroyed in many areas. During the winter, the voles invaded barns and farm buildings, destroying all kinds of stored food. In 1932-35 voles and lemmings were in epidemic proportions in over 104,000 square kilometres (40,000 square miles) of what was then the southern U.S.S.R. This resulted in extensive damage to wheat and other field crops and to orchards.

Rats will eat almost anything that humans eat. Perhaps the most serious damage is to the seeds of grain plants, both before and after harvesting. Grain stored on farms is often not only eaten by rats but also rendered unsuitable



Figure 30: Range of body-plan variation of larger Rodentia.

for human consumption by being mixed with rat droppings. Food that has reached warehouses in cities is also eaten by rats, and here the excess of damage over the amount actually consumed is even greater. It is estimated that rats damage about twice as much grain as they eat. About 23 kilograms (50 pounds) of grain are required to support a rat for a year, so the total cost is about 70 kilograms (150 pounds) per rat per year. Rats sometimes demonstrate a fondness for animal food and have been known to kill several hundred baby chickens in a single night. Eggs are frequently eaten, and even full grown hens, baby pigs, and lambs may be killed. Rats will also follow a farmer in the planting season and dig up newly planted seeds.

Depredation by rats

Rats also do extensive damage in searching for food or in making nests. They gnaw holes to gain entrance to barns, warehouses, or houses, or through walls once they are inside. Clothing, upholstered furniture, and other textiles are gnawed to provide nest-building materials. In areas of cities where extensive quantities of garbage and other refuse are available, the rat populations may exceed the human populations, and rat infestation is one of the commonest complaints of slum dwellers. Such rats are very apt to bite sleeping humans, especially children, and fatal attacks on babies have been known to occur.

The house mouse has food preferences similar to those of the Norway rat but is not as abundant nor as large. Other rodents are much less likely to attack stored supplies of human food, but they do eat considerable quantities before harvesting. This is true of such forms as voles, field mice, and squirrels. The seeds of both wild and domesticated grasses are a major food source for large numbers of wild rodents. Because of the quantities of these that are eaten, such animals are often considered a major threat by grazing interests, but the benefit that rodents give by collecting seeds into underground store houses, where they may later sprout and grow into new plants, certainly comes close to balancing any damage they cause to man's interest.

The role of rats in the "Black Death"

Transmission of diseases. *Plague.* Rodents serve as reservoirs for a number of diseases that may be transmitted to humans by arthropod agents. The most devastating of these is bubonic plague. This disease is fundamentally a disease of rodents, especially rats, transmitted from one rodent to another by an intermediate host, the flea, which also serves to transmit the disease to humans. An epidemic (called a "pandemic" because of the totality of infection in the human population) that seems to have been plague spread over Europe in the 6th century AD. If this epidemic was indeed plague, it must have involved an unusual rodent host because it antedated the arrival of *Rattus* in Europe. The epidemic known as the Black Death originated in Mesopotamia about the middle of the 11th century, and spread to Europe, particularly in the 14th century, being accompanied by the spread of rats throughout the Continent. It has been estimated that 25,000,000 people died of that pandemic of the plague in Europe. The latest pandemic originated in southwestern China in the late 19th century and was spread all over the world by the rat populations of oceangoing ships. Plague is controlled by controlling rat populations, particularly those on ships, and preventing shipborne rats from reaching land. Although plague epidemics have been brought under control, there still remain numerous foci of infection.

The plague bacillus (*Pasteurella pestis*) can infect a variety of other rodents, and foci have been established among native rodents other than *Rattus* in many parts of the world. This disease is referred to as "sylvatic plague" to distinguish it from the basically urban occurrence of ratborne plague. Over 80 species of ground-living and burrowing rodents are known to be involved in sylvatic plague, including ground squirrels, some cricetids (e.g., voles, lemmings, muskrats), several murids (e.g., Old World mice), and a few others, including guinea pigs.

Tularemia. Tularemia is primarily a disease of lagomorphs (rabbits and hares) and secondarily of rodents. It has been reported from all parts of the United States, and seems to have spread from there to many other parts of the world. Among rodents, it is carried by ground squir-

rels, tree squirrels, prairie dogs, chipmunks, muskrats, and beavers, and is transmitted largely by ticks.

Rickettsial diseases. Three rickettsial diseases—murine typhus, Rocky Mountain spotted fever, and tsutsugamushi disease—involve rodents as reservoirs. Murine typhus, which is much less severe than epidemic typhus, is transmitted to man by fleas, primarily from *Rattus*. It is probably almost universal in tropical and subtropical areas. Rocky Mountain spotted fever, which apparently originated in the northwestern United States, has spread over much of that country, and very similar if not identical diseases are found in Mexico, South America, and Africa. It is transmitted to humans by tick bites, and is presumably endemic in a considerable number of rodents as well as other animals. The cottontail rabbit is believed to be a major reservoir animal. Tsutsugamushi disease, or scrub typhus, is transmitted to humans by the bite of a rat mite. It is found in the East Indies and Southeast Asia. The reservoir seems to be primarily rodents of the genus *Rattus*. A variety of other diseases of lesser importance are transmitted to humans from a reservoir in rodents, largely *Rattus*, by rat bites, ticks, or in other ways.

Damage by burrowing and gnawing. Rodents have been accused, in many parts of the world, of damage to agricultural interests because of their burrows. Cattlemen in the western United States supported campaigns to poison prairie dogs (*Cynomys*) because horses or cattle might (and occasionally did) break their legs in prairie-dog holes. But perhaps the most striking case of damage from rodent burrows has been caused by the American muskrat, introduced into European fur farms, from which it escaped as early as 1905, spreading widely over much of Europe, from Great Britain and Brittany to the Ukraine. It burrows in the banks of streams, and has caused damage especially by making burrows in the banks of drainage ditches, canals, and in dikes.

The ever-growing incisors of rodents need to be used regularly on hard substances to wear them off. This need results in the gnawing of lead pipes or telephone transmission cables, or the making of holes in boxes, walls, and stored inedible items.

Benefits derived from rodents. *Trapping for furs.* The beaver (*Castor canadensis*) was extremely important in the development of the western United States and Canada, the trappers of the first two-thirds of the 19th century being primarily interested in beaver skins. Most of the initial exploration of the West was performed by beaver trappers. The beaver skin became the basic unit of currency over much of the western United States. The trapping was so extensive and effective that beavers were exterminated over much of the area, but they are making a slow recovery under current conditions of protection. The fur of the chinchilla (*Chinchilla laniger*) was of such value that these rodents were almost exterminated in Argentina. Muskrat and nutria have also been extensively hunted for their fur.

Rats in the laboratory. Among the best laboratory mammals, for all types of biological, medical, and psychological investigation and for testing of new drugs, are the laboratory rats and mice, normally albino strains of wild species of the Norway rat and the house mouse. Guinea pigs (*Cavia cobaya*) are also widely used. There are a number of other laboratory rodents of lesser importance, of which the Syrian or golden hamster (*Mesocricetus auratus*) is the most widely used. Some strains of rats are nearly unique among nonhuman animals in being susceptible to dental caries, and most of the experimental work on tooth decay has been performed on them.

Rodents as valuable furbearers

ECOLOGY AND NICHE RELATIONSHIPS

Habitat and locomotion. The most typical members of the order are the small, ground-living rodents, of the type generally called rats or mice, similar to what was probably the ancestral condition for the order. While these generally stay on the ground, most can climb shrubs and bushes with ease and many climb trees; among the smallest rodents, harvest mice climb stalks of wheat to reach the seeds on which they feed. This central generalized adaptive type of rodent, the scamperer with limited burrowing or arboreal adaptations, includes the family Muridae (Old World rats

The basic scampering rodent

and mice, now worldwide), the family Cricetidae (field mice, wood rats, voles, lemmings, muskrats, to name a few; found in most parts of the world other than Australia), some Sciuridae (the ground squirrels and chipmunks), the nonleaping members of the New World family Heteromyidae (pocket mice); some South American rodents (degus, spiny rats, and chinchilla rats), and the South African rock rats (Petromuridae). Although there are many adaptive variants among the scampering rodents, there is a strong tendency for them to inhabit either open country or limited areas of woodland. Rodents of this type are found over essentially the entire land surface of the world. In most cases individual or family territories are set up, and intruders are successfully kept at bay by a variety of threatening and defiant postures. Normally there is no actual combat. The territories are, of course, of varying size, related to the available food supplies and the size of the animals. Individual ground squirrels have been observed over a territory of 2,500 square metres (about 0.6 acre). Changes with time in the territories of chipmunks (*Eutamias*) indicate rather striking seasonal changes, with an increase in the area toward the end of the foraging season (perhaps to allow greater accumulation of food), and very little shift of the territory from one year to another. With few exceptions, these terrestrial rodents make burrows near the centre of the territory where they may spend as much as half of the time, where the young are born, and where the adults hibernate, in forms where hibernation occurs. Normally most young are produced than there is room for, and as a result there is extensive migration of subadult or young adult individuals. For most rodents, the lengths of these migrations are unknown, but it has been shown that muskrats will migrate as much as 30 kilometres (19 miles), including several kilometres cross-country, over a period of one to two weeks, looking for suitable new stream or swamp sites.

From occasional tree climbing it is a relatively short evolutionary step to an arboreal habitat, such as that occupied by the tree squirrels or, with an increase in size and decrease in rapidity of movement, by the New World porcupines. Tree squirrels are, essentially, scamperers that have developed the capability of holding on to bark with their claws. They are as much at home on the ground as in trees and show no special anatomic modifications for tree life. Their nests may be either in hollow trees or other sheltered places (e.g., barns) or built of twigs and leaves among the smaller branches of trees. In Europe and adjacent parts of North Africa and Asia, the dormice (Gliridae) have similar habitat but spend an even larger part of their lives in the trees. Many murids and cricetids (rats and mice, in the broad sense) are about as fully arboreal in their habits as are the tree squirrels.

Scampering arboreal animals are faced with the necessity, from time to time, of moving rapidly from one tree to another and do so by leaping. It is generally agreed that gliding forms have evolved their skin fold (patagium), which assists them in gliding, as a direct adaptation to this type of habitat. Patagia have evolved, probably several times independently, among the Sciuridae (squirrels) of both the Old World (*Pteromys* and others) and North America (*Glaucomys*) and in the quite distinct and unrelated African family of scaly-tailed squirrels (Anomaluridae). The patagium of the anomalurids is supported at the anterior end by a long, slender cartilage, which has apparently evolved from the olecranon process of the ulna (the posterior bone of the forearm) at the elbow.

Larger, heavier arboreal forms such as the New World porcupines move slowly along the trunks and main branches of trees. One of these, in South and Central America, has developed a prehensile (grasping) tail that assists in climbing.

A large proportion of the small, ground-living rodents construct underground nests, where the young are born and reared. Often there are adjacent chambers for food storage. In many cases two or more exits provide an opportunity for escape from predators. From ancestors such as this, many lines of rodents have evolved into habitual burrowers that make extensive subterranean tunnels and rarely or never come to the surface, foraging, feeding, mat-

ing, and raising their families underground. Among the numerous varieties of burrowing rodents are the prairie dogs, colonies of which were formerly abundant in the western United States. The mouth of the burrow is surrounded by a volcano-shaped mound of dirt, on which a prairie dog frequently sits upright, alert for potential enemies. The animal forages on the ground surface, returning to the mound or to the burrow to eat. Voles (*Microtus*) may make lengthy runways 50 to 75 millimetres (two or three inches) below the surface, arching up the ground surface over them. The pocket gophers (*Geomys* and *Thomomys*) are even more effective burrowers, throwing up mounds 30 centimetres (12 inches) or more across and spending most of their life beneath the ground. They use the incisors as well as the front legs in digging. During times of heavy snow, they may burrow through snow, feeding on plants above ground. Similar types of burrowing occur in the Old World mole rats (Spalacidae). An extreme form of burrowing is found in the African mole rats or blesmols (Bathyergidae), in which the incisors extend forward, with a fold of skin closing the mouth behind them, permitting their use as the primary digging tools, having replaced the forelimbs in this respect. The burrows of a single blesmol form a crisscross of tunnels, at several levels below the ground surface, that may occupy an area of a thousand or more square metres. All these burrowing forms have short, heavy limbs used in digging, short tails, and small eyes. In one genus of blesmol, *Heliophobius*, the eyelids have grown together so that the eyes do not form images. These animals remain permanently below ground. Many burrowing rodents are solitary, but a number of them live in colonies, unlike other members of the order.

Relatively few rodents have become large animals. Several of the South American caviomorphs have done so, however, and have evolved as quadrupedal terrestrial cursorial animals. Usually these have the claws modified into or toward hooves, reducing the number of toes on each foot to four or occasionally three and tending to make each foot symmetrical about the median plane. Among these are the capybara (the largest living rodent), nutria, paca (*Cuniculus paca*), and agouti. The guinea pig is a smaller relative with the same adaptations. Such animals normally escape predators by running away, although both capybara and nutria will take to the water. The Old World porcupines are slow-moving terrestrial quadrupeds that are protected, as are the New World porcupines, by the quills, which are elongate, pointed hairs.

Many rodents have become bipedal jumpers, enabling them to escape their enemies with long leaps. These forms include the kangaroo rats (*Dipodomys* and *Microdipodops*) of western North America, the jumping mice (*Zapus* and *Napaeozapus*) of eastern North America and eastern Asia, the jerboas (Dipodidae) and gerbils (*Gerbillus*) of North Africa and Southeast Asia, the saltatorial dormouse (*Selevinia*) of Kazakhstan, the chinchilla of the southern Andes, and the springhaas or Cape jumping hare (*Pedetes*) of South Africa. All of these are small, except *Pedetes*, which stands 30 centimetres (12 inches) tall. Most of the saltatorial rodents have very highly enlarged tympanic bullae—the bony covering of the middle ear—which is not true of other rodents. It is probable that this characteristic is an adaptation to life in a desert, where an amplification of the noise made by approaching predators may be of considerable importance. Certainly the large middle-ear cavity greatly increases the amplification of sounds. The fact that the bullae are particularly large in leaping forms suggests further that there may be some relationship between the size of the bulla and the balancing and stabilizing activities of the ears.

Although a number of rodents are at least partially aquatic, special adaptations are rare, being found only in the beavers. Beavers have large hind legs, with strong webs between the toes and an enlarged first digit, as long as the other toes, making very efficient paddles. Swimming is carried out primarily by the feet, but the broad flat tail may be used on occasion for sculling. A number of other rodents that are semiaquatic have webbed feet and soft, dense underfur but few other aquatic adaptations. These semiaquatic forms include muskrats, water voles,

Giant rodents

Gliding rodents

Burrowing rodents

capybaras, nutrias, South American fish-eating rats, and the water rats of Australia and New Guinea.

Seed-eating rodents

Food habits. As a whole rodents are herbivorous, although most of them will eat animal food on occasion. The diet includes a wide variety of plant foods, although seeds are the favourite item. Many are particularly fond of the seeds of the various grasses, both wild and domestic. These may be eaten on the spot, but large quantities are usually carried home in the cheek pouches that are common among rodents. Seeds carried home may be eaten in the safety of the burrow or stored for later use. Such storage is an important item to the plants concerned, as it places the seeds in a favourable location for sprouting, if the rodent either forgets about them or himself becomes food for a predator. Many rodents, such as squirrels, are fond of nuts and can open even the hardest black walnuts to get at the seed within. Acorns, seeds, fruits, berries, young leaves, buds and shoots, green leaves, tubers, and bulbs are eaten by a wide variety of rodents. Burrowing rodents tend to concentrate on the roots, tubers, and bulbs that they encounter in making their tunnels, although some burrowers forage on the surface as well. The larger caviomorphs, such as the capybara, paca, mara (*Dolichotis patagona*), and agouti (*Dasyprocta aguti*), are essentially grazing animals, eating all kinds of green vegetation, young stems, and fruits. Marmots (*Marmota*) have rather similar food habits. Old World porcupines are essentially omnivorous, including carrion in their food. New World porcupines eat leaves, twigs, buds, and bark. Beavers cut aspens, willows, and poplars, using leaves, buds, and bark for food; they also eat a variety of aquatic plants. Many desert-living rodents have a diet that, in the dry season, may be almost exclusively seeds, from which they are able to obtain enough water for survival, using the water produced by metabolism. African cane rats eat coarse grasses and shrubs, one of their favourite foods being sugarcane. A number of rodents, especially during hard winters, strip the bark from trees. Because of their large numbers, rodents can seriously damage the natural plant cover.

On the other hand, many rodents will, on occasion, eat animal food. Muskrats eat freshwater mussels and crayfish, as well as a wide variety of dead or dying water animals. Many rodents, will, at least occasionally, eat birds' eggs, nestlings, or insects. Several genera of South American cricetids live mostly on aquatic animals, of which fish make up the largest part. The grasshopper mice (*Onychomys*) of western United States and adjacent parts of Mexico and Canada are almost exclusively carnivores. Their food consists primarily of insects and earthworms, but may include birds and even other rodents that they kill. In the Arctic portion of their range, ground squirrels eat any carrion available, including their own dead relatives or even walrus or whale meat. They particularly like meat with a high fat content, presumably for the energy value.

Longevity. Rodents are generally considered to have a very short life span. This is certainly true in nature, where rodents are one of the primary food sources of carnivorous birds and mammals. Particularly among the small, mouselike rodents, the average life expectation at birth must be only a few weeks or months; relatively few individuals live much more than a year; and an animal two years old has reached a ripe old age. It is clear that this limitation is largely imposed on the animals by the activities of predators. Cricetids, for example, which normally live less than two years in the wild, have been kept up to five years in captivity, chipmunks and ground squirrels up to seven or eight years, gerbils and dormice over five years, guinea pigs about eight years, chinchillas and North American porcupines about 10 years, and woodchucks, pacas, and the Old World porcupine even longer.

Predation on rodents. Because rodents are one of the most important items in their food supply, carnivorous birds (especially hawks and owls) and mammals provide a very strong selective force acting on small rodents. Many raptorial birds use their eyes for hunting and are able to detect any rodent slightly less well concealed than the average. It has been shown experimentally that owls selectively capture light-coloured mice over dark ones on a dark background and dark over light on a light back-

ground down to a light level of 10^{-7} of a footcandle. As a result of such selective pressure, rodents have evolved fur colours that closely match the environmental background where they live. For example, a late Pleistocene flow of dark lava in New Mexico is now inhabited by melanistic (dark) rats and mice, in strong contrast to the normally coloured rodents surrounding it and to the almost white rodents found in the nearby White Sands area.

Among mammalian carnivores, all but the very largest normally eat rodents, at least as an important part of their diet, and many will start teaching their young how to hunt at the expense of small rodents. The abundance of rodents, moreover, affects the numbers of their predators in at least some instances. The English ecologist Charles Elton has shown that there are striking parallels between cyclic abundance of lemmings and snowy owls and foxes in northern Labrador, the lemming maxima and minima preceding those of the predators by about six months to a year.

BEHAVIOUR

Gnawing. One of the most general habits of rodents is gnawing, much of which serves merely to wear down the incisors. These teeth grow continuously, at rates that have been found to range upward from two millimetres (0.08 inch) per week in the few species in which rates have been measured. Gnawing is performed by all rodents at frequent intervals. It involves the fore-and-aft movement of the lower jaw, the upper incisors holding a hard object and the lower incisors cutting against it. The same movements may be used with nothing held in the teeth, the lower incisors merely alternating between cutting against the rear side of the upper incisors and shifting forward so that the uppers cut against the rear of the lowers. Because the hard enamel is limited to the anterior side of the teeth, such use wears away the softer dentine at a faster rate, leaving the enamel to form a sharp chisel edge. Rodents prefer to use their incisors on hard objects, which provide part of the abrasion. These may be nuts, bark, tree trunks (as in beavers), often bones (which may also be gnawed for a supply of calcium), or even human possessions such as boards in houses and barns or even metal telephone cables.

Function of gnawing

Nesting and denning. *Food storage.* A large proportion of rodents build underground homes, with a central nest chamber, in which they sleep, raise their young, and, often, hibernate. Many others have similar nests in grass or trees. Special food storage areas often are associated with such living quarters; for burrowing forms, these food areas are special chambers off the main passageways. Ground squirrels, kangaroo rats, and other rodents bring grass seeds back to such chambers in their cheek pouches (on the inside of the mouth of ground squirrels; outside the mouth and fur lined in the Heteromyidae and Geomyidae), and will often store several times as many seeds as they will eat during the season when they rely on their stores. Excavation of one kangaroo rat den, from which a single rat weighing 149 grams (about five ounces) was removed, showed nine underground storage chambers, containing from one to nine litres (one to eight quarts) of seeds, with a total of almost 39 litres (35 quarts). Kangaroo rats, in addition, dig storage pits, about 2.5 centimetres (one inch) in diameter and the same in depth, near their burrows, which they fill with seeds. In one case, 875 such caches were located in an area of 5.1 square metres (55 square feet) adjacent to a single den. It has been estimated that seed collection in this manner may result in the loss of as much as a quarter of the grain crop in some parts of the world, and at times of rodent plagues, the destruction locally may approach totality. On the other hand, such underground stores are an important source of germinating seeds for wild grasses in steppe regions such as western North America and Central Asia.

The habit of squirrels of carrying acorns and nuts to hollow trees or barns, or digging holes in the ground in which the nuts are placed, is well known. The animal's memory of the location of these stores is poor, but if enough nuts or acorns are hidden, the animal will be able to find a sufficient number to keep him alive during the winter.

Limits of life span

The North American pack rat or trade rat (*Neotoma cinerea*) is attracted by bright and shining objects, which it picks up to carry home to its nest, a jumble of sticks, twigs, grass, and assorted collectors items. A popular superstition is that this animal is a fair businessman, who, on seeing something he wants, always leaves a replacement that is, in his opinion at least, of equal value. The fact is that, while carrying one trophy, the rat may see another that is more attractive and so puts down the first to pick up the second, leading the human who has lost a possession to believe that the rat was carrying out a "business deal."

Constructions by beavers

Beaver dams and lodges. Among the best known activities of rodents are those of beavers, which go to great lengths to store food, cutting alder and other food trees into 0.5- to 2.5-metre (1.6- to 8-foot) lengths and floating the logs in their ponds or embedding one end of these sticks in the mud of the ponds formed by their dams. These sticks make an available food supply during the winter when the ponds are frozen. In order to provide an aquatic habitat, safe from their normal predators, beavers build dams that often run several hundred feet in length, and as much as two metres (6.6 feet) high on the downstream side. The largest dam reported, from near Three Forks, Montana, was 652 metres (2,140 feet) long, 4.3 metres (14 feet) high at the highest point, and seven metres (23 feet) thick at the base. The dams are made of twigs, branches, and logs from which the bark has been eaten, piled on the ground with alternating layers of gravel or mud. When finished, the upper surface is plastered with mud, usually dug from the bottom of the pond above the dam, to make it watertight. The pond behind the dam serves as the refuge and first defense of the beaver. Consequently, beavers work continually to maintain the dams, by repairing, raising, and lengthening them, a process that may extend over several generations. Abandoned beaver ponds gradually silt up, forming very characteristic meadows.

Beavers also build lodges in the ponds behind their dams or burrows in the banks. Lodges are built up in the same manner as dams, and they are enlarged along with the rising water level of the ponds as the dam is raised. For ventilation, however, the upper part of the lodge is much more loosely constructed than the dam. The floor is kept above the water level, and there may be two levels, the lower for feeding, the upper for sleeping. The entrance to the lodge is through a tunnel, the outer end of which is below water level. Bank lodges are burrows in the bank, with the inside arranged like that of lodges. Again, access is gained through a tunnel with the mouth below water level. Burrows may extend more than nine metres into the bank.

Beaver canals are less well known but highly remarkable. They are dug to permit the transport of birch or aspen logs for food supplies from the sources where they are cut to the ponds. The logs are floated in the water, where the beaver is much safer than on land. Such canals are dug whenever they can be made with little or no damming. They vary from 0.3 to 1.2 metres (one to four feet) in width, up to 0.6 metres (two feet) deep, and canals as long as 223 metres (732 feet) have been reported.

Beavers live on the bark of trees, that of aspen being their favourite, but birch, cottonwood, willow, and alder also are used. When these are in short supply other plants are eaten, but conifers are rarely cut except for building materials. In one night a single adult beaver can fell a seven- to 10-centimetre (three- to four-inch) poplar, cut it into sections one to three metres long, and drag the logs to the water. Normally, beavers do not cut up logs over 15 centimetres (six inches) in diameter, although they will fell trees considerably larger. The largest tree on record cut by beavers was 115 centimetres (46 inches) in diameter. The cutting of trees is apparently done at random, with all gnawing likely to be on one side, unless the tree has a large diameter, in which case it may be gnawed all the way around. Because trees normally fall downhill, hence toward a pond, this is probably the origin of the legend that beavers can make trees fall as they wish.

Hibernation. Many rodents use their nests or burrows for shelter during unfavourable seasons, especially winter. There is considerable variation, even among closely

related animals, in the amount of activity that occurs at such times and how often animals wake and eat some of the stored food. In a considerable number, including chipmunks, pocket mice (*Perognathus*), some jerboas, and some field mice, there is extensive torpor in wintertime, the animals sleeping extensively, but waking two or three times a day to eat. Hamsters (*Cricetus* and *Mesocricetus*) are deep hibernators, going into a sound sleep from which they awaken only periodically to eat. The most pronounced type of hibernation involves the accumulation of extensive amounts of body fat during the late summer and early fall, and deep hibernation with few or no awakenings during the winter. This condition occurs in marmots, some ground squirrels, a large proportion of dormice (whose English name is derived from the French *dormir* meaning "to sleep"), the jerboa (*Allactaga*), and most or all zopodids (jumping mice and birch mice). Estivation, or summer sleep, is rare in rodents but occurs in woodchucks and some ground squirrels. No suggestions of hibernation occur in murids, caviomorphs, or phiomorphs, any other tropical rodents, or rodents with a recent tropical or subtropical ancestry.

Population movements. Several types of rodents, subject to periodic fluctuations in numbers, reduce the amount of overpopulation by migration. Generally, the migration is simply the movement of juveniles away from a home area that has become fully populated. Migration of such animals may be quite extensive, but, because it is usually an individual phenomenon, it is often not noticed. Such animals are subject to very high mortality from predators because they are moving in unaccustomed territory and without the shelter of a permanent home. Muskrats and beavers are examples that are more readily noted, because they move from one suitable stream to another and may travel as much as several miles across country looking for a suitable home. There are many reports of major migrations. In August 1946, gray and fox squirrels migrated out of an area in northwest Wisconsin, apparently due to a shortage of acorns. When they reached streams and rivers, they crossed them by swimming, many drowning in the wider rivers. In 1935 a mass movement of gray squirrels was noted from the area east of the Hudson River, to the west. Hundreds (one report says 2,000) of drowned squirrels were found along the west bank of the Hudson River.

The most famous migrations of rodents, however, are those of lemmings, especially in Scandinavia. There tend to be two migrations per year, one in spring and one in autumn. The exact causes of the movement are not known, but a change of habitat with season seems to be the most important single factor. It has been believed that the migrations are caused by overpopulation and a shortage of food, or by claustrophobia (a feeling of mental anguish from the crowded conditions), but these probably are not valid causes.

The lemming migrations are individual activities rather than unified movements, as is popularly supposed. Careful observations of migrations show that each individual acts alone, but that there may be groups of several individuals within a few minutes, followed by a gap of 10 minutes or more. In autumn, at least, the migrations normally occur at night. Individual lemmings generally move fairly rapidly in a generally constant direction, although they will occasionally stop to feed or groom. Rates of migration have been calculated as about one metre per second (2.25 miles per hour), but such speeds are not sustained, and it is probable that distances of eight-10 kilometres (five-six miles) per day are normal. The migrations tend to be oriented by roads or paths made by humans, trails made by reindeer, or other similar routes, although the general direction seems to be outward in all directions from the general source area and to be determined before the movement begins. Contrary to the older reports and popular belief, lemmings do not rush forward without hesitation into water when reaching a shoreline. They apparently try to avoid swimming, if possible, running up and down the shore looking for a narrow crossing or an area where they can cross on ice. Apparently, they enter water willingly only if they can see the silhouette of the far shore. This is in marked contrast to activities of their relative, the water

Lemming migrations

Types of dormancy

vole (*Arvicola*), which always heads for the water when threatened. When swimming lemmings use their hind legs almost exclusively but swim very high out of the water. Lemmings become exhausted and drown after swimming 15 to 25 minutes in water with waves about 15 centimetres (six inches) high. Observations along the Swedish-Finnish border show that the lemmings had no problem crossing a strait 200 metres (650 feet) wide on a calm night, but that considerable numbers drowned on a windy night. A lemming was observed swimming at approximately the middle of a lake two kilometres (1.2 miles) wide. Reports of lemmings swimming out into the North Sea and being found alive 16 kilometres (10 miles) or more from shore must be viewed skeptically.

Recent studies in Scandinavia indicate that the migrations are a normal part of the activity of the lemmings. They serve not only as a means of dispersal of the species but also to permit the animals to occupy different habitats at different seasons.

Rodents are often thought of as defenseless animals, but the incisors are very sharp and powerful, are activated by powerful jaw muscles, and can inflict a deep wound. Among themselves, rodents fight for territories, the possessor of a territory usually being able to drive interlopers away. They also fight very effectively in self-defense, as indicated by the popular expression, "fighting like a cornered rat." A full-grown wharf rat is capable of defending itself against a determined tomcat. Even animals as small as lemmings are able to discourage attack by weasels.

REPRODUCTION

Flexibility
in breeding
habits

The reproductive habits of rodents are exceedingly varied in nature and are capable of being even further modified in domesticated or partially domesticated forms. Many, especially the larger types, reproduce once a year. Others produce several litters during a single season. Some have only one or two young at a time; others, large numbers. Some are capable of reproduction at very early ages, others (particularly members of the Caviomorpha) not until after a considerable period of growth. Most are born naked and helpless, and are cared for in nests; a few are able to run and keep up with their mothers almost at once. Most rodents are polygamous; some mate for the duration of a single breeding season; and a few (beavers, for example) have permanent mates.

Beavers do not breed until the second January or February after birth, producing a litter of two or four young after about 12 weeks of gestation. The capybara has a litter that is, on the average, slightly larger, after 15 to 18 weeks gestation; the viscacha, a smaller animal, has two young in the spring about 22 weeks after breeding. Nutrias breed twice a year after reaching an age of one year, if the winter temperature does not fall below 16° C (60° F) for lengthy periods, and produce litters averaging about five young. The European porcupine (*Hystrix cristata*) breeds in the spring, with one to four young being born about 16 weeks later, and some of its tropical relatives (other members of the Hystricidae) produce two litters a year. The North American porcupine (*Erethizon dorsatum*) of the caviomorph family Erethizontidae, has a single young after 30 weeks gestation. Most of the smaller members of the order, however, have considerably higher reproductive rates. In many cases there may be more than one litter

per year, with high numbers per litter, examples being shown in Table 2.

This rather high rate of breeding is intensified by the fact that, in many of the smaller rodents, sexual maturity is reached at an early age, normally earlier in the females than in the males. The females breed when less than a year old in many squirrels, some pocket gophers, and the pack rat. In most of the Cricetidae, however, the young appear to reach sexual maturity considerably earlier—harvest mice (*Reithrodontomys*) in five weeks; the hamster (*Cricetus*) in six weeks; and voles (*Microtus*) in six to seven weeks.

In the United States, wild populations of the house mouse reproduce throughout the year, with an average of 5.5 litters and 31 young per female per year in buildings, and 10.2 litters and 57 young per year on farms. In the laboratory mouse, the mean litter size has been reported as ranging from 4.5 to 7.4, depending on the strain. The second litter is the largest, after which there is a steady decrease. Litter size has been reported to have a range of two to 12, and as many as 19 healthy embryos have been removed just before term from a single female. There may be five or more litters per year. The gestation period is about 20 days, and the first mating usually occurs at seven to 10 weeks of age, though it has been reported that, in one strain of albino mice, the first estrus occurs, on the average, at 39 days and results in about 50 percent pregnancies. The breeding span of most mice in the laboratory is 12 to 18 months. Lactating females may become pregnant, but the gestation will be lengthened by one to two weeks.

Wild Norway rats breed throughout the year, taking advantage of the sheltered environment available from living in association with humans. The size of the litter is correlated with the size of the mother. In the laboratory female rats reach sexual maturity in 33 to 120 days. The number of young varies with the strain, two typical ones having litters averaging 6.5 and 8.9. Normally, a female's second or third litter is the largest, and the litter size decreases rapidly after the tenth. Over a number of generations of laboratory life in a colony of gray rats, however, the total number of young per female grew from an average of 23 to an average of 63.

The laboratory guinea pig has a gestation period of about 68 days, a very long period for such a small animal but characteristic of caviomorphs. There are usually three or four young per litter. The age of puberty varies from 55 to 70 days.

Although most rodents seem to increase their reproductive potential under domestication, this is not true of chinchillas. In the wild, they produce litters of one to six after about 15 weeks pregnancy; domesticated ones seem normally to have only one offspring at a time.

FORM AND FUNCTION

In most aspects rodents are relatively primitive, not highly specialized mammals. The skeleton is usually that of a quadrupedal, scampering mammal, not far from the basic placental stock. The digits are clawed and usually five in number, with little or no opposability of the first digit. The limbs are not far from the same length, although the hind legs are always somewhat longer. The animals are adapted to a broadly herbivorous diet. The main specialization of the digestive system is the possession of a large cecum (a blind pouch) at the junction of the large and small intestines.

Specializations for gnawing. As opposed to this basic primitiveness, the rodents have developed, as their basic innovation, a gnawing mechanism of an efficiency that is not approached by that of any other mammals. The glenoid fossa, the concavity where the lower jaw articulates with the skull, slopes downward, from rear to front, and there are no processes (as there are in all other mammals) at the anterior and posterior limits of the fossa that restrict the anteroposterior movement of the jaw. The jaw, therefore, is capable of functioning when pulled to the rear, separating the incisors from each other, and permitting chewing or grinding by the cheek teeth. The jaw can also be pulled forward and downward, separating the

Reproductive
potential

Table 2: Reproductive Rates of Certain Rodents

animals	litters/year	young/litter
North American gray squirrel (<i>Sciurus carolinensis</i>)	2	2-3
Eurasian gray squirrel (<i>S. vulgaris</i>)	2	5-7
Red squirrels (<i>Tamiasciurus</i>)	2	4-6
Pocket gophers (<i>Geomys</i>)	1-2	1-4
Pocket mice (<i>Perognathus</i>)	1-2	2-8
Kangaroo rats (<i>Dipodomys</i>)	1-3	2-4
Deer mice (<i>Peromyscus</i>)	up to 4	2-7
Hamsters (<i>Cricetus</i>)	several	6-12
Lemmings (<i>Lemmus</i>)	several	3-9
Muskrats (<i>Ondatra</i>)	2-5	5-7
Voiles (<i>Microtus</i>)	up to 13	4-8

Source: E.P. Walker, *Mammals of the World* (1964) and S.A. Asdell, *Patterns of Mammalian Reproduction* (1964).

Evolutionary shifts in jaw muscle attachments

cheek teeth and bringing the tips of the incisors together for gnawing.

These shifts in the position of the jaw, as well as the functioning of the jaw in either position, are brought about by a combination of the actions of the various jaw muscles, of which the temporal, pterygoid, and masseter are the most important. The forward and backward (anteroposterior) movements, including those of gnawing, were caused, primitively, by all three muscles, all of which

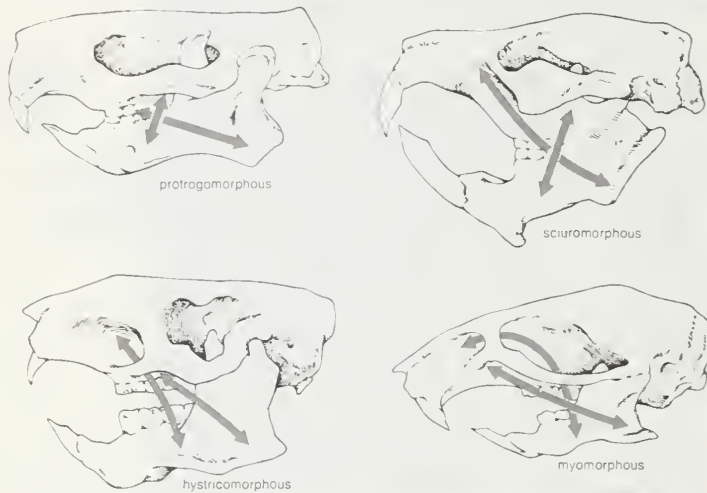


Figure 31: Skulls of the four basic types of rodents. Arrows show positions of branches of the masseter muscle.

have a slight anteroposterior component. But a shift of the anterior muscle, the masseter, began to occur during the Eocene Epoch (about 50,000,000 years ago). Originally, in a condition called protrugomorphous, the masseter ran from the zygomatic arch (cheekbone) to the lower jaw. Parts of the muscle, retaining the original points of insertion on the lower jaw, shifted their origins forward onto the snout, doubling the length of the muscle and greatly changing its direction of pull. In some rodents it was the lateral portion of the masseter that shifted forward onto the face, compressing the infraorbital foramen (an opening in the bone), through which pass nerves and blood vessels, between the muscle and the snout (the sciurumorphous condition); in others it was the deeper portion of the masseter that shifted forward, through the infraorbital foramen (presumably following an initial enlargement of the foramen), onto the face, resulting in further enlargement of the foramen until it may now be larger than the eye socket (the hystricomorphous condition); and in still others, a combination of the two shifts has occurred (the myomorphous condition). Although these changes are of importance in determining rodent taxonomic relationships, they are probably not as controlling as was once thought, as each structural condition may have originated more than once.

Continuous growth of incisors

But the main structural change from the primitive mammalian condition has occurred in the incisors. In all known rodents, beginning with the very earliest, the incisors are reduced to a single pair in each of the upper and lower jaws. These teeth grow throughout life, the enamel cap being restricted to the anterior 30 to 60 percent of the perimeter of the tooth. The rate of incisor growth has been measured in several rodents, and varies from about two to three millimetres per week (0.08–0.12 inch) in nonburrowing forms to five millimetres (0.2) per week in pocket gophers, in which the incisors are used for digging as well as for gnawing. In hibernating rodents, the incisors continue to grow, but at a greatly reduced rate: If the incisors do not meet, due to deformity or breakage, the teeth will continue to grow indefinitely. Since the lower teeth are arcs of large circles, the lower incisors simply curve forward and upward, becoming completely nonfunctional. But the upper incisors are large arcs of much smaller circles. With growth, they may spiral outward, like the spirals on a handle-bar moustache, or they may grow around and upward through the throat between the lower jaws, growing through the snout and eventually locking the jaws

closed. Either condition is fatal, but rodents have survived in nature for a long time with such malformed incisors.

Because of the continual growth of their incisors, rodents have an unusually heavy need for calcium, which, together with a need for abrasion of the incisors, explains the frequency with which rodents gnaw bones.

Adaptations for specialized locomotion. Many rodents have shifted their locomotion from the characteristic scampering type. In burrowing rodents, the limbs are short and massive, and the hands, particularly, are widened and strengthened with heavy claws. Many, but not all, burrowing rodents use their incisors to help in digging. In such cases, a fold of skin is likely to close the mouth behind the incisors, and the radius of curvature of both upper and lower incisors is increased, so that they extend forward from the mouth, cutting the soil an appreciable distance in front of the animal. Eyes and ears are generally reduced in size in these forms.

There is little difference, structurally, between the terrestrial scampering rodents and the arboreal scamperers, such as squirrels. The claws of the latter are sharper, but locomotion is basically the same. A few rodents, such as the New World porcupine, actually climb (as would a human) rather than running up and down the bark, as does a squirrel. Only one rodent, the tree porcupine (*Coendou*) of Central America and northern South America, has a truly prehensile tail (although several rodents have semiprehensile tails that they can wrap around branches for extra support). A number of gliding rodents, "flying squirrels" of one sort or another, some of which are not closely related to the true squirrels (*Sciuridae*), are supported in the air by a fur-covered membrane (patagium) extending between the fore and hind limbs and usually between the hind limbs and the tail.

A rather common development among rodents has been the reduction of the forelimbs and an increase in the size of the hind limbs and tail for locomotion by means of a series of kangaroo-like leaps, which in several forms can exceed 4.5 metres (15 feet).

Adaptations for water conservation. With their wide geographic distribution, rodents have invaded desert regions in all parts of the world. Many of these have developed the ability to function on limited supplies of water. This has been achieved in several ways. Some are nocturnal, remaining asleep during the heat of the day in burrows deep enough to reach ground moisture, thus reducing evaporation. Others seek out succulent desert plants for food. But it has been shown that a number of desert rodents (kangaroo rats and jerboas, for example) have the ability to live exclusively on seeds, with no external water source, and to obtain all of their needed water from their metabolism. Such forms also produce a greatly reduced quantity of extremely concentrated urine and a small amount of feces. Two species of spiny mouse (*Acomys*) were investigated in Israel. Both inhabit the same desert areas, but one is diurnal and the other is nocturnal. Both produce urine the urea content of which is extremely high. Neither species can survive on a diet of barley without water to drink, although gerbils and jerboas from the same area showed no ill effects after four weeks without water. The diurnal spiny mice, however, were able to survive and grow with no water except seawater, whereas the nocturnal ones could not tolerate such a salt concentration.

EVOLUTION AND PALEONTOLOGY

Rodents are relatively poorly represented in collections of fossils, in spite of their great abundance at the present time. They have often been overlooked because of their small size, but modern intensive exploration for fossils usually results in a much more abundant representation of small mammals, especially rodents.

Paleocene. The earliest known rodents come from the late Paleocene (about 57,000,000 years ago) of North America, by which time they had already acquired all of the diagnostic features of the order. The ancestral family, the Paramyidae, was also present in Europe, where it first appeared in the earliest Eocene. There was very rapid diversification of the order during the Eocene, initially involving the Paramyidae, but with other families soon

Para-
myidae,
the earliest
rodents

appearing. Most of the Eocene rodents were protrogomorphous, the masseter muscle restricted to its primitive origin on the cheekbone. By the middle Eocene of Europe, rodents with advanced types of masseters were present, and others occurred in the late Eocene of North America. The skeleton of the paramyids was basically that of a generalized scampering animal, approximating a rat in its method of locomotion. Before the end of the Eocene it is probable that both leaping and burrowing variations had arisen, although there is some uncertainty on this point.

Oligocene. A major gap in the knowledge of rodent evolution occurs at the Eocene-Oligocene boundary (about 38,000,000 years ago). A number of modern families with advanced types of jaw muscles appeared at about the same time in North America, Europe, or Asia, including the Cricetidae, Heteromyidae, Geomyidae, Castoridae, and Ctenodactylidae, with no good indications at present as to their geographic or ancestral sources. (The Ctenodactylidae are known from the Eocene of Kazakhstan, however.) At about the same time, rodent groups appeared suddenly in South America and Africa (Caviomorpha and Phiomorpha, respectively). Both groups are hystricomorphous and hystricognathous and have often been considered to have been related, though others consider them to have acquired their similarities independently.

Miocene. The majority of the living families of rodents had appeared by the Oligocene, and most of the remaining ones occur in the Miocene (around 26,000,000 years ago), so that the later part of the evolution of the rodents was largely a diversification of the various groups. The ancestors of a wide variety of specialized modern forms can be recognized from teeth, but most of them are not known from skeletons, and the evolutionary status of the modern locomotor conditions is not well known. Highly evolved burrowing rodents were certainly present in the Oligocene in Mongolia, however, and less specialized ones occurred in North America. Ancestral kangaroo rats had attained a leaping ability by the Miocene of North America. Gliding and burrowing forms are known from the Miocene of Africa. In South America, the differentiation of the descendants of the invaders present in the early Oligocene proceeded rapidly, and the Miocene forms occupied nearly all potential rodent niches on that continent.



Figure 32: Reconstruction of *Epigaulus*, a primitive Pliocene rodent.

Appearance of Muridae

Later evolution of rodents. The latest rodents to differentiate apparently were the members of the family Muridae. This group appears suddenly in early Pliocene deposits of Europe, probably having invaded that continent from Asia, the southeastern part of which, together with the adjacent East Indies, is now considered the evolutionary centre of the family. The Muridae expanded rapidly all over the Old World. An isolated incisor is known from the Pliocene (approximately 4,000,000 years ago) of New Guinea, where considerable local evolution has occurred. The murids reached Australia late in the Cenozoic Era (they are unknown there before the Pleistocene), and underwent a rapid evolution, developing arboreal, burrowing, and leaping forms, paralleling the evolution of the entire order in the rest of the world.

When North and South America were united in the beginning of the Pleistocene (about 2,500,000 years ago), there was a major invasion of South America by cricetids,

which expanded rapidly and became highly diverse in that area.

An interesting aspect of rodent evolution are the very rapid diversifications (evolutionary explosions) that occurred whenever highly specialized members of the order reached areas that had been previously uninhabited by gnawing mammals (South America and Africa in the Oligocene, New Guinea in the Pliocene, Australia in the Pleistocene). These, apparently, resulted in a rapid spread into most of the major available ecologic niches in a matter of, at most, 2,000,000 or 3,000,000 years. The differentiation of the South American cricetines since the beginning of the Pleistocene has not been quite as fast, though still explosive, presumably because there already were numerous rodents on that continent. The initial diversification of the rodents during the Eocene seems to have taken place at a considerably slower rate, probably because they had not yet become as thoroughly adapted for gnawing as was the case later.

Because of their small size and great abundance, rodent fossils are becoming progressively more important as guide fossils that enable accurate separation of successive geologic horizons.

CLASSIFICATION

Distinguishing taxonomic features. Many features have been used to determine the relationships of the subgroups of the rodents. The most generally used are the differentiations of the jaw muscles or the modifications of the skull and jaws resulting from muscular changes. The masseter muscle may have its primitive position, arising from the cheekbone (protrogomorphous); its lateral division may have shifted forward onto the snout, compressing the infraorbital foramen against the snout (sciurormorphous); the medial division may have moved forward through the infraorbital foramen, onto the snout, greatly enlarging the foramen (hystricomorphous); or both branches of the masseter may have shifted (myomorphous). In most rodents the angular process at the rear of the lower jaw extends downward from the ventral side of the alveolus of the incisor (sciurognathous); in some, all of which are also hystricomorphous, the angle arises from the side of the alveolus (hystricognathous). Other features that have been given important weight in delimiting the major subgroups include: crown patterns of the cheek teeth (premolars and molars); histologic structure of incisor enamel, which may have one layer (uniserial), a few (pauciserial), or many (multiserial); structure of the penis or of the baculum (os penis); structure of the male genital tract, including particularly the presence or absence of an outgrowth, the sacculus urethralis; fusion or lack thereof of two of the ear ossicles, the malleus and incus; and the pattern of development of the extra-embryonic fetal membranes.

ANNOTATED CLASSIFICATION

The classification presented here is based upon that of American paleontologist A.E. Wood. The rodents are one of the outstanding examples of a group in which closely similar structures have evolved numerous times independently, a phenomenon known as parallelism, with the result that there is very strong disagreement among taxonomists working with the group as to how the order should be subdivided.

Groups indicated by a dagger (†) are known only from fossil remains. Dental formulas indicate the number of pairs of teeth present in the upper jaw (above the line) and lower jaw (below), the figures representing, respectively, pairs of incisors, canines, premolars, and molars. When a partial formula is given the type of tooth is indicated by I, C, P, or M.

ORDER RODENTIA

Gnawing mammals with ever-growing incisors, reduced to a single pair in both upper and lower jaws and with the enamel limited to the anterior face of the incisors, so that wear maintains a sharp chisel edge; glenoid fossa slanted, with neither preglenoid nor postglenoid process, masseter the principal jaw muscle; dental formula reduced, never exceeding $\frac{1 \cdot 0 \cdot 2 \cdot 3}{1 \cdot 0 \cdot 1 \cdot 3} = 22$. About 350 living genera and 2,400 living species; over 400 extinct genera have been described.

Suborder Sciuromorpha

Masseter primitively either limited to zygoma, or rarely with only a slight forward displacement, but in the squirrels shifted forward onto face (sciuromorphous). Always sciurognathous. No sacculus urethralis. Malleus and incus bones of ear separate. Incisor enamel pauciserial or uniserial. Paleocene to Recent.

†*Family Paramyidae*

Paleocene to early Miocene; Northern Hemisphere. Cheek teeth cuspidate, clearly derived from basic mammalian tribosphenic type. Tympanic bullae became co-ossified with the skull several times independently, within the family. Incisor enamel pauciserial, with one possible exception. Presumably ancestral to the rest of the order. One subfamily incipiently hystricognathous. From size of a mouse to as large as a beaver. Locomotion, when known, scampering, with possible incipient leaping in one form.

†*Family Sciuravidae*

Eocene of North America; one genus from Eocene of Central Asia. Cheek teeth formed of 4 transverse crests rather than of separate cusps; locomotion probably scampering; mouse-size to rat-size. Incisor enamel pauciserial. Perhaps ancestral to several mouse- or ratlike groups.

†*Family Ischyromyidae*

Late Eocene to Oligocene of North America. Masseter muscle has begun to move forward onto the snout in some forms, in primitive position in others. Cheek teeth 4-crested. Incisor enamel uniserial. Locomotion scampering. Size of a gray squirrel or somewhat larger.

†*Family Cylindrodontidae*

Middle Eocene to Oligocene of North America. Oligocene of Central Asia. Burrowing rodents with high-crowned to ever-growing cheek teeth, based on a pattern of 4 transverse crests. The most specialized forms used the incisors for digging, as shown by forward extension of these teeth. Incisor enamel uniserial. Size from that of a chipmunk to that of a marmot. Presumably derived from North American Sciuravidae.

†*Family Protoptychidae*

Late Eocene of North America. Skull with highly inflated auditory bullae, suggesting, by analogy with other rodents, that these were leaping animals. Cheek teeth high crowned, with pattern of 4 crests. Size of a gray squirrel.

Family Aplodontidae (mountain beaver, or sewellel, and fossil relatives)

Late Eocene to Recent of North America. Oligocene and Miocene of Europe, Pliocene of Asia. Single living species restricted to wet areas from British Columbia south to San Francisco Bay. Extreme hypsodont cheek teeth, based on cuspidate rather than crested pattern. Incisor enamel uniserial. Burrowers. Head and body length of living species (*Aplodontia rufa*) 300 to 460 mm (11.8–18.1 in.), with a short tail; weight 900 to 1,800 g (2–4 lb).

†*Family Mylagaulidae*

Miocene to Pliocene of North America. Strongly hypsodont cheek teeth, with greatly enlarged premolars, whose continual growth forced 1 or more of the anterior molars out of the jaws. Incisor enamel uniserial. Some individuals had paired horns on the snout. Whether these are taxonomic or sexual characters has not yet been demonstrated. Powerful burrowers. About the same size as *Aplodontia*.

Family Sciuridae (squirrels, chipmunks, and marmots)

Oligocene to Recent of Northern Hemisphere, Miocene to Recent of Africa, Recent of South America. Highly developed sciuromorphous pattern, with masseter muscle extending forward along the side of the snout and compressing the infraorbital foramen. Postorbital process of frontal separates the eye from the temporal muscle. Cheek teeth typically very similar to those of the Paramyidae, but some forms have rather complicated tooth patterns. Incisor enamel uniserial. About 70 living genera.

Suborder Myomorpha

Lateral branch of masseter muscle usually displaced forward alongside of snout, forcing infraorbital foramen against side of snout; in most forms, deep masseter penetrates a variable distance through upper part of infraorbital foramen. Sciurognathous. Incisor enamel uniserial. Locomotion scampering, jumping, arboreal scampering, or fossorial. A few forms are partially adapted to an aquatic life. Malleus and incus never fused. No sacculus urethralis. Cheek teeth usually reduced in number, with only a single fossil genus known that retains the primitive rodent dental formula; except for the Geomyoidea, the premolars are either lost or greatly reduced; the last molars are occasionally lost as well. Most members of the group are

small, from the size of a mouse to that of a rat (head and body length ranges from 50 to 300 mm (2–11.8 in.), except in the Rhizomyidae, which may be somewhat larger).

Family Cricetidae (field mice, deer mice, voles, lemmings, muskrats)

Early Oligocene to Recent of Europe and North America, middle Oligocene to Recent of Asia, late Pliocene to Recent of South America, Pleistocene of Madagascar, Recent of Africa. Deep masseter expanded through upper part of infraorbital foramen to face; no premolars, molars $\frac{3}{3}$; cheek-tooth pattern based on cusps arranged into 5 transverse crests in both upper and lower teeth; teeth low crowned to very high crowned. Scampering, fossorial, arboreal scampering, jumping, or partly aquatic.

Family Muridae (Old World rats and mice)

Early Pliocene to Recent of Europe and Asia, late Pliocene to Recent of East Indies, Pleistocene to Recent of Africa and Australia. Introduced throughout the world by man. Deep masseter as in Cricetidae; cheek teeth normally the 3 molars, but one subfamily, the Hydromurinae of Australia and New Guinea, has lost the third molars, reducing the cheek teeth to M_2^2 ; tooth pattern based on rounded cusps, arranged in transverse rows, but derivable by modification of a pattern like that of primitive cricetids; low to medium high crowned teeth. Scampering, arboreal scampering, occasionally fossorial or semiaquatic.

Family Heteromyidae (pocket mice, kangaroo rats and mice)

Oligocene to Recent of North America. Recent of northern South America. These and the next two families are sciuro-morphous, with no penetration of the infraorbital foramen by the masseter, and seem to be closely related. Fur-lined cheek pouches, opening beside the mouth and reaching back to the shoulders, are used for the transportation of food to underground storage areas near the nests. Central stock of family scampering; specialized ones leap with their hind legs. Generally small, mouselike; length of head and body from 55 to 180 mm (2.2–7.1 in.). Cheek teeth, which include $P_1^1 M_2^2$, are low crowned to very high crowned. Pattern of teeth based on 2 transverse rows of 3 cusps each. Derivable from late Eocene members of the Eomyidae.

Family Geomyidae (pocket gophers)

Early Miocene to Recent of North America. Similar to heteromyids in regard to jaw muscles, dental formula, and tooth pattern, and in possession of fur-lined cheek pouches. Highly adapted burrowing animals. Cheek teeth have already become high crowned in the Miocene, and those of the living forms grow throughout the animal's life, with the enamel reduced to a plate on the anterior side of the upper teeth and the posterior side of the lowers.

†*Family Eomyidae*

Late Eocene to late Pliocene of North America; late Eocene to Pleistocene of Europe. Sciuromorphous jaw muscles; cranial anatomy and jaw structure very similar to heteromyids. Cheek teeth usually low crowned, with pattern of 5 crests, rather similar to that of the cricetids, although some members of the family developed teeth suggesting that they might be ancestral to heteromyids. Cheek teeth usually $P_1^1 M_2^2$, but one genus had P_1^1 . Skeleton and habits unknown.

Family Zapodidae (birch and jumping mice)

Late Oligocene to Recent of Eurasia, early Miocene to Recent of North America. A late Eocene genus from California is often (perhaps incorrectly) placed here. The Old World birch mice are arboreal scampering forms, as were probably all the earlier fossil members of the family. The jumping mice are as saltatorial as the kangaroo rats. Infraorbital foramen of medium size, with the deep masseter passing through it but with the superficial masseter remaining on the zygomatic arch. Cheek teeth $P_1^1 M_2^2$, with pattern very similar to that of cricetids, with which the fossils have frequently been confused.

Family Dipodidae (jerboas)

Late Oligocene to Recent of Asia, late Miocene to Recent of Europe, Pleistocene to Recent of North Africa. Masseter muscle like that of zapodids but with larger deep division. Highly saltatorial rodents of the steppes and deserts of the Old World, with extremely inflated auditory bullae. Cheek teeth $P_1^1 M_2^2$, high crowned in all living and most fossil forms. Tooth pattern suggestive of that of zapodids, from which the family is probably descended. General tendency to elongate the hind foot and to fuse the 3 median metatarsals to form a cannon bone, as part of the leaping adaptation.

Family Spalacidae (mole rats)

Early Pliocene to Recent of Europe. Pleistocene to Recent of Western Asia and North Africa. Highly specialized burrowing rodents with short, powerful legs and no external tail. Eyelids permanently closed. Masseter muscle a short distance forward

on face. Cheek teeth reduced to M_3^1 , (or possibly $P_1^1 M_2^1$). Ancestry unknown.

Family Rhizomyidae (African mole rats, bamboo rats)

Oligocene of Europe, early Miocene to Recent of Asia, and Pleistocene to Recent of Africa. Myomorphous, with considerable forward extension of the masseter. Burrowing animals with powerful legs and short tails, superficially resembling pocket gophers (Geomyidae). Cheek teeth high crowned to ever-growing, consisting of $P_{1-0}^{1-0} M_3^3$. The European Oligocene *Rhizospalax* has features suggesting that this family and the Spalacidae are related. Otherwise, relationships unknown.

Suborder Caviomorpha

Deep branch to masseter has shifted its origin forward onto the face, passing through the very large infraorbital foramen. Angle of jaw hystricognathous. Malleus and incus fused; usually a sacculus urethralis in male genital tract. Incisor enamel multiseriate. Cheek teeth usually $P_1^1 M_2^1$, but occasionally the milk premolar is retained throughout life, the formula being dP_1^1, M_2^1 . Early Oligocene to Recent of South America, Pleistocene to Recent of North America and West Indies.

Family Octodontidae (octodonts, degus)

Early Oligocene to Recent of South America, Pleistocene to Recent of the West Indies. Small for caviomorphs, generally somewhat ratlike in appearance. Scampering to fossorial in habits. Occur from sea level to elevations of more than 3,000 m (10,000 ft) in the southern half of South America. The teeth range from low crowned to ever-growing, the enamel of the grinding surface arranged either in the shape of a kidney or figure 8. The burrowing forms use the incisors in digging. The genus *Platypittamys* from the early Oligocene of Patagonia is very close to being a common ancestor of all the Caviomorpha.

Family Echimyidae (spiny rats)

Oligocene to Recent of South America. Pleistocene to Recent of West Indies and Central America. Ratlike in appearance. Usually with spiny fur, although there are a few exceptions. Cheek teeth rooted, with a pattern formed by transverse folds of enamel. In all but the earliest known members of the family, the milk or deciduous premolars are retained throughout life, and permanent premolars never make an appearance. Inhabit moist forested regions; climbing or scampering; one genus is burrowing.

Family Ctenomyidae (tuco-tucos)

Pliocene to Recent of South America. The single living genus inhabits all of the southern half of South America, from sea level to elevations over 4,000 m (13,000 ft). They are fossorial and presumably derived from octodontids. Body form and size are very similar to those of the North American pocket gophers, with powerful digging muscles in the forelimbs and long, powerful claws.

Family Abrocomidae (chinchilla rats or abrocomes)

The common name comes from the soft underfur and ratlike appearance. Pliocene to Recent of South America. The living species inhabit mountainous areas of Peru, Bolivia, Argentina, and Chile. Presumably evolved from octodontids. The cheek teeth grow throughout life.

Family Chinchillidae (chinchillas and viscachas)

Early Oligocene to Recent of South America. One living genus (*Lagostomus*) lives in lowlands of Argentina; the other two (*Lagidium* and *Chinchilla*) live at elevations of about 800 to 6,500 m in the southern half of the Andes. The fur of *Chinchilla* is very soft and dense—that of the other genera less so. All gregarious, but especially the lowland viscachas (*Lagostomus*), which occupy colonies resembling those of prairie dogs. The lowland viscachas are large, head and body length being 470 to 660 mm (18.5–26 in.); the mountain genera are 230 to 320 mm (9–12.6 in.; *Lagidium*) and 225 to 380 mm (8.9–15 in.; *Chinchilla*). Cheek teeth ever-growing, but details of pattern lost early in life.

Family Capromyidae (hutias, coypus)

Middle Pliocene to Recent of South America, Pleistocene to Recent of West Indies, introduced in southern United States and parts of Europe. Differ from the Chinchillidae, from which they are probably derived, in having high-crowned to ever-growing cheek teeth, but with the details of the pattern persistent. The hutias of the West Indies resemble large rats. The coypus of South America are considerably larger, looking rather like a large muskrat. The fur is excellent, and the flesh is widely eaten in South America.

Family Dasyproctidae (pacas and agoutis)

Early Oligocene to Recent of South America, Pleistocene to Recent of West Indies, Recent of Central America. Large rodents (head and body length 320 to 800 mm [12.6–31.5 in.]), with weight up to 10 kg (22 lb). Limbs modified for cursorial locomotion, with the lateral toes, especially on the hind foot, reduced in size. Flesh is very palatable. Pacas in moist forested areas from Mexico to southern Brazil; agoutis in the same regions (plus the Lesser Antilles), but not restricted to forested areas.

rial locomotion, with the lateral toes, especially on the hind foot, reduced in size. Flesh is very palatable. Pacas in moist forested areas from Mexico to southern Brazil; agoutis in the same regions (plus the Lesser Antilles), but not restricted to forested areas.

Family Dinomyidae (pacaranas)

Early Miocene to Recent of South America; Pleistocene of West Indies. There is only a single living genus inhabiting the lower elevations of the northern half of the Andes. The animal is among the largest of living rodents. A considerable number of fossil forms, once thought to represent a separate family (Heptaxodontidae) probably belong here.

†*Family Elasmodontomyidae*

Pleistocene to Recent of Puerto Rico and northern Lesser Antilles. These forms are all extinct, but survived until about the time the area was colonized by the Spaniards. Medium-sized ground-living rodents, with ever-growing cheek teeth formed of a series of enamel plates inclined at about 45° to the long axis of the jaws.

†*Family Eocardiidae*

Early Oligocene to middle Miocene of South America. Cheek teeth medium to high crowned. Size about that of a guinea pig.

Family Caviidae (guinea pigs, cavies, and maras)

Late Miocene to Recent of South America. Cheek teeth ever-growing, with a simplified crown pattern of alternating V's. Digits reduced to 4 on the front and 3 on the hind foot. Cavies with short legs, ears and tails; maras superficially resembling hares.

Family Hydrochoeridae (capybaras)

Early Pliocene to Recent of South America, Pleistocene to Recent of Central America, Pleistocene of West Indies and southern United States. Most members of this family are extinct, there being but a single living genus, with 1 species in eastern Panama and the other in South America east of the Andes and north of the Rio Paraná. Forest dwellers, capable of escape by running, but tend to retreat to water when closely pursued. Derived from Caviidae and have the same reduction of digits. Largest living rodents, head and body length reaching 1.2 m (3.9 ft); they may weigh over 45 kg (100 lb).

Family Erethizontidae (New World porcupines)

Early Oligocene to Recent of South America, Pleistocene to Recent of North America, recent of Central America. Large, slow-moving, heavy-bodied rodents, with some hairs modified into sharp, barbed spines that are easily detached when they make contact with an enemy. Cheek teeth rooted, with a simple pattern of reentrant folds, essentially unchanged since the Oligocene. One South American genus has a prehensile tail.

Suborder Phiomorpha

A group of families of African rodents the ancestry of which can be traced back to the early Oligocene of Egypt. Most of them are hystricomorphous and hystricognathous, and they are often considered very closely related to the Caviomorpha. Multiserial incisor enamel; malleus and incus fused and sacculus urethralis present in living forms. Cheek teeth $\frac{4}{1}$, consisting (in all but a few of the earliest types) of the last deciduous premolar, retained throughout life, and M_1^1 .

†*Family Phiomyidae*

Oligocene to Miocene of Africa. These rodents reached Africa when it was isolated from the rest of the world and differentiated to become highly diverse, including scampering, burrowing, and perhaps arboreal and leaping forms. In the Miocene other rodents reached Africa, and only a few lines of phiomyids were able to survive the resulting competition. Several Oligocene genera retained permanent premolars, although in 1, at least, they apparently never were functional.

Family Petromuridae (rock rats or dassic rats)

Pleistocene to Recent of southern South Africa. Ratlike in body form, the single species inhabits rocky hills, living in narrow crevices in rocks.

Family Thryonomyidae (cane rats)

Miocene to Recent of Africa and early Pliocene of India. Body form and habitat generally similar to a muskrat, inhabiting marshes and the borders of streams and lakes. There is only a single living genus, spread over Africa south of the Sahara. Head and body 350 to 610 mm.

Family Bathyergidae (blesmols or African mole rats)

Miocene to Recent of Africa. The most completely fossorial of all rodents, with short, powerful limbs and heavy claws. The incisors are highly procumbent, and in some forms are used in burrowing. In some genera the growing base of the upper incisors has shifted backward to a level behind the rear of the upper cheek teeth. Only slight penetration of in-

fraorbital foramen by masseter (probably secondary reduction); angle hystricognathous but peculiar. One genus (*Heterocephalus*) has the pelage reduced to scattered short hairs. External ears are greatly reduced; eyes are reduced and the eyelids usually kept closed. The number of cheek teeth reaches $\frac{5}{6}$; it is by no means clear which teeth are involved. These animals spend essentially all their lives in their burrows; some are colonial.

Families of uncertain relationships

The remaining 10 families of rodents are of completely uncertain relationships, and are best treated as individual families.

Family Hystricidae (Old World porcupines)

Pliocene to Recent of Asia, Africa, and Europe, Pleistocene to Recent of East Indies. Spined forms, resembling the New World porcupines in appearance, but less arboreal in habits; quills without barbs. The ease with which even a slight touch loosens the quills, leaving them inserted in the attacker, is the origin of the belief that these animals are capable of shooting their quills like arrows, as reported, for example, by Marco Polo. Hystricognathous and hystricomorphous; malleus and incus fused; incisor enamel multiserial; sacculus urethralis present; cheek teeth consist of P_1^+ , M_3^+ , all exhibiting a crown pattern broken up into numerous small cusps. The origin and relationships of the group are not clear, but they may have originated in southern Asia.

Family Castoridae (beavers)

Early Oligocene to Recent of Europe and North America, late Oligocene to Recent of Asia. Sciuro-morphous and sciurognathous rodents, with cheek teeth (P_1^{2-1} , M_3^+) progressively becoming high crowned and ever-growing. Restricted at present to a single genus (*Castor*), formerly widespread in Eurasia and North America, but greatly reduced as a result of very intensive trapping for the very fine, soft fur. Aquatic animals, with webbed feet and broad, flattened tail. Numerous fossils from the Tertiary do not seem to have been aquatic but rather to have been fossorial. In western Nebraska an early Miocene beaver seems to have constructed spiral burrows ("Devils' corkscrews").

†*Family Eutypomyidae*

Rodents from the Oligocene of North America sometimes (but probably erroneously) thought to have been related to the beavers, but with very complexly folded enamel on the cheek teeth, and with a peculiar foot structure of uncertain function (very slender inner digits and heavy outer ones).

Family Anomaluridae (scaly-tailed "squirrels")

Miocene to Recent of Africa. Sciurognathous but hystricomorphous rodents; arboreal, and most living forms with a membrane extending from front to hind legs to tail, used in gliding from one tree to another. A long cartilaginous support, derived from the ulna in the region of the elbow, supports the anterior end of the membrane. Head and body length about 60 to 430 mm (2.4–16.9 in.).

Family Ctenodactylidae (gundis)

Oligocene of Mongolia, Miocene of India, and Miocene to Recent of Africa. Hystricomorphous and sciurognathous rodents; malleus and incus fused; perhaps a sacculus urethralis; incisor enamel multiserial. Quite diverse in Oligocene of Mongolia, but the fossils do not suggest close relationship to any possible ancestors. Scamperers, with body form similar to that of a guinea pig; 4 toes on each foot. Head and body length 160 to 240 mm (6.3–9.4 in.).

Family Pedetidae (springhaas or Cape jumping hare)

Miocene to Recent of Africa. Hystricomorphous and sciurognathous; malleus and incus not fused; incisor enamel multiserial. The single living genus lives in eastern and southern Africa. The Cape jumping hare is a large rodent, with elongate hind limbs, jumping like a kangaroo. The cheek teeth (P_1^+ , M_3^+) are ever-growing, with the pattern preserved only in completely unworn teeth. Head and body length about 350 to 430 mm (13.8–16.9 in.).

Family Gliridae (dormice)

Middle Eocene to Recent of Europe, Recent of Asia, Miocene to Recent of northern Africa. Sciuro-morphous and sciurognathous; uniserial incisor enamel. Small arboreal rodents, similar in appearance to the smaller squirrels. Generally have partial hibernation. It has recently been demonstrated that they can be traced back to an ancestry among small European paramyids of early to middle Eocene; as a result, although they are similar to sciurids in being sciuro-morphous, they are clearly of independent origin. Head and body length about 60 to 190 mm (2.4–7.5 in.).

Family Selevniidae (jumping dormice)

Probably related to dormice; a single genus recently discovered in the deserts of Central Asia. Elongate hind legs and a long tail. Head and body length 72 to 96 mm (2.8–3.8 in.).

†*Family Pseudosciuridae*

European, middle Eocene to Oligocene. Hystricomorphous and sciurognathous; incisor enamel pauciserial. Cheek teeth low crowned. They can be derived from European early Eocene paramyids, and gave rise to the Theridomyidae. The gradual development of the hystricomorphous condition can be traced in this group.

†*Family Theridomyidae*

Abundant European late Eocene and Oligocene rodents, with low-crowned to ever-growing cheek teeth. Incisor enamel pauciserial or uniserial. They have had a central position in many theories of the origin and relationships of the various hystricomorphous groups (Caviomorpha, Phiomorpha, Hystricidae, Anomaluridae, Ctenodactylidae, and Pedetidae), but there is no evidence to show that they were actually related to any of these.

Critical appraisal. The rodents are one of the most clearly demarcated of all mammalian orders. No animals are known, either living or fossil, where there is question as to whether or not they belong to this order. Within the order, however, there is great disagreement as to the interrelationships of various families or groups of families. There is no question that the Caviomorpha, as listed above, are related, and there seems to be little question that the Phiomorpha are a natural group. There is disagreement as to whether the two suborders should be united with each other, or with the Hystricidae. If the three are united, the combination should be called the suborder Hystricomorpha (the oldest name for such a group). The Anomaluridae, Ctenodactylidae, Pedetidae, Pseudosciuridae, and Theridomyidae are included in the Hystricomorpha by some authors. The relationships of all five families are very uncertain (except that pseudosciurids clearly gave rise to theridomyids). Much uncertainty exists as to whether the families here included in the Myomorpha are a natural unit. There is considerable diversity within the group. Many authors unite the Muridae and Cricetidae into a single family; still others separate the microtines from the cricetids as a third family. The two families are recognized here largely as a matter of convenience—about 100 genera each of murids and cricetids are known; combining them makes an unwieldy family including over half the order. The microtines were quite obviously derived in late Pliocene times from cricetids; if they are recognized as a distinct family, it is the only family of either animals or plants known to have originated so recently.

Currently, attempts are still being made to find satisfactory bases on which to subdivide the order. The incompleteness of knowledge of fossil rodents is a major handicap that will certainly be overcome in the future. The classification presented here is an attempt to give what seems to be a reasonable interpretation of the present knowledge of the subject. (A.E.W.)

Carnivora (cats, dogs, bears, skunks, seals, walruses)

The order Carnivora includes 10 families of living mammals: Canidae (dogs, wolves, jackals, and foxes), Ursidae (bears), Procyonidae (raccoons), Mustelidae (skunks, mink, weasels, badgers, and otters), Viverridae (civets and mongooses), Hyaenidae (hyenas), Felidae (cats), Otariidae (eared seals), Odobenidae (walrus), and Phocidae (earless seals). The term carnivore is frequently applied by mammalogists to members of this order and is employed in that sense in the present article. In a more general sense, a carnivore is any animal (or even, occasionally, a plant) that eats the flesh of other animals, as opposed to a "herbivore," which eats plants. Although the Carnivora are basically meat eaters, a substantial number, especially among bears and procyonids, feed extensively on vegetable material.

GENERAL FEATURES

Importance of carnivora. There is probably no other group of animals more familiar to man than mammals belonging to the order Carnivora. The more popular do-

Uncertain relationships

mesticated pets of man, the dog and the cat, are both derived from wild members of this order. The majority of luxurious natural furs (ermine, mink, sable, otter) worn by women of fashion, as well as the furs of utility used by primitive man, come from members of the Carnivora. Many of the animals that attract the largest crowds at circuses and zoos are members of this order. Most of the large, dangerous carnivores are the objects of hunters, who wish to obtain a trophy grizzly bear, polar bear, or tiger. The man stranded in the wild fears the wolf, dhole, bear, tiger, leopard, or jaguar, virtually the only predators capable of injuring him. The stockman is concerned about possible depredations upon his herd by nature's large predators.

The carnivores, the meat eaters, form the apex of the pyramid of life, the food pyramid, and thus are basic to the existence of a "balance of nature" in a world unmolested by man. In today's world, man's world, this precarious balance was first upset by the extermination of many carnivores. Formerly considered undesirable because of their predation on game animals, carnivores are now recognized as necessary elements in natural systems. Far from depressing game numbers, carnivores improve the stability of game populations by keeping numbers within the carrying capacity of the food supply. As a result, individual animals are better fed and less subject to disease. Many of these predators dig dens and provide burrows in which other forms of wildlife can take refuge. Those carnivores best known for their burrow building are the badgers and the skunks. Digging results in the mixing of soils and the reduction of the runoff of water during rains.

Intelligence and training

Carnivores rank high on the scale of intelligence among mammals. The large size of the brain, compared to that of the animal, is an indication of their superior mental powers. For this reason, these animals are among the easiest to train for entertainment purposes, as pets, or as hunting companions. Their highly developed sense of smell supplements the sharper vision of man. Dogs are the most common carnivores trained for hunting, but the cheetah, caracal, and ferret have also been used to some extent. In China the otter is trained to drive fish under a large net, which is then dropped and pulled in. Carnivores, dependent for survival upon their ability to prey upon living animals in a variety of situations, have evolved a relatively high degree of learning ability (see **LEARNING, ANIMAL**).

Carnivorous mammals tend to establish territories, areas defended against others of the same species; herbivorous ones, which eat vegetation, are less apt to do so. Territories are often exclusive, defended by the residents against other animals of their own kind. Such areas may sometimes be marked by secretions produced by anal or scent glands.

There is a wide range of social patterns among carnivores. Some (bears, raccoon, red and gray foxes, genets, most cats, and most mustelids) are strictly solitary, except during the breeding season. Some remain paired throughout the year (black-backed jackal and lesser panda) or occasionally hunt in pairs (gray fox, crab-eating fox, and kinkajou). Other carnivores, such as the wolf, African hunting dog, dhole, and coatimundi, normally hunt in packs or bands. Still others form sedentary colonies during the breeding season (sea lions, fur seals, and elephant seals), during a somewhat larger part of the year (sea otters), or all year round (meerkats).

Form and function. The smallest living member of the Carnivora is the least weasel (*Mustela nivalis*), which weighs about two ounces. The largest terrestrial form is the Alaskan grizzly (*Ursus arctos*), weighing up to 780 kilograms (1700 pounds). The largest aquatic form is the elephant seal (*Mirounga leonina*), which may weigh 3,640 kilograms (four tons). Most carnivores weigh between four and eight kilograms (nine and 18 pounds).

Most members of the order are terrestrial. Some, such as the sea otter, river otter, and polar bear, spend most of their lives in water. The pinnipeds, or seals, are more aquatic than other members of the order. Aquatic or semi-aquatic forms tend to have body specializations such as a streamlined body and webbed feet for this mode of life.

Dentition

Carnivores, like other mammals, have a number of different kinds of teeth: incisors in front, followed by canines,

premolars, and molars in the rear. Most carnivores, especially those that feed exclusively on meat, have carnassial, or shearing, teeth that function in slicing meat and cutting tough sinews. The carnassials are usually formed by the fourth upper premolar and the first lower molar, working one against the other with a scissorlike action. Cats, hyenas, and weasels, all highly carnivorous, have well-developed carnassials; while the bear and procyonids, which tend to be omnivorous (eating both plants and animals), and the seals, which eat fish or marine invertebrates, have little or no modification of these teeth for shearing. The teeth behind the carnassials tend to be lost or reduced in size in highly carnivorous species. Most members of the order have six prominent incisors on both the upper and lower jaw, two canines on each jaw, six to eight premolars, and four molars above and four to six molars below. Incisors are adapted for nipping off flesh. The outer most incisors are usually larger than the inner ones. The strong canines are usually large, pointed, and adapted to aid in the stabbing of prey. The premolars always have sharply pointed cusps, and in some forms (*e.g.*, seals) all the cheek teeth (premolars and molars) have this shape. Except for the carnassials, molars tend to be flat teeth utilized for crushing. Terrestrial carnivores that depend largely on meat, such as weasels, cats, and hyenas, tend to have fewer teeth (30-34), the flat molars having been lost. Omnivorous carnivores, such as raccoon and bear, have more teeth (40-42). Seals have fewer teeth than terrestrial carnivores. In addition, seals exhibit little stability in the numbers of teeth; for example, a walrus may have from 18 to 24 teeth.

Several features of the skull are characteristic of the order Carnivora. The articulating surfaces (condyles) on the lower jaw are transverse, their axis at right angles to that of the head, forming a half-cylindrical hinge that allows the jaw to move only in a vertical plane but with considerable strength. The clavicles (collarbones) are either reduced or absent entirely and, if present, are usually embedded in muscles without articulation with other bones. This allows for a greater flexibility in the shoulder area and prevents breakage of the clavicles when the animal springs on its prey.

The brain is large in relation to the weight of the body and contains complex convolutions characteristic of highly intelligent animals. The stomach is simple, and a blind pouch (cecum) attached to the intestine is usually reduced or absent. Since animal tissues are in general simpler to digest than plant tissues, the carnivore's dependence on a diet with a high proportion of meat has led to less-complex compartmentalization of the stomach and a decrease in the length and folding (surface area) of the intestine. The teats are located on the abdomen along two primitive lines (milk ridges), a characteristic of mammals that lie down when nursing.

Distribution and abundance. Carnivores are found worldwide. Terrestrial forms are absent from most oceanic islands, though the coastlines are usually visited by seals (pinnipeds). Except for the dingo (which probably was introduced by aboriginal man), Australia has no native terrestrial members of the Carnivora. Man has taken his pets, as well as a number of wild species, to most islands. For example, a large population of the red fox now inhabits Australia, having been introduced there by fox hunters.

Since carnivores are large and depend on meat, there must be fewer carnivores in the environment than those animals that form their diet. The maintenance of established territories limits the number of predators to the carrying capacity of prey populations. In general, carnivores have a population density of approximately 0.4 per square kilometre (one per square mile). By comparison, omnivorous mammals average about eight per square kilometre (20 per square mile), and herbivorous rodents attain densities of up to 40,000 per square kilometre (100,000 per square mile) at peak population. Existing at relatively low densities, carnivores are vulnerable to prey, population fluctuations, habitat disturbance, and predation by man. In some cases the mobility and adaptability of carnivores has enabled species to shift ecological roles and survive the changes brought about by human activities; in other instances less-flexible species have become extinct.

World carnivore population

The world population of aquatic carnivores (pinnipeds) has been estimated to be between 13,000,000 and 27,000,000 animals. Estimates of the total numbers of terrestrial carnivores have not been made. The yearly sale of furs indicates a harvest of at least 6,000,000 pelts of wild carnivores; this is equalled only by the number of mink raised and harvested on fur farms (approximately 8,000,000 annually in the mid-1960s). There are probably around 50,000,000 to 60,000,000 terrestrial members of the Carnivora present today. Some species (many cats, pandas, and bear, and some seals) are becoming quite rare and near extinction.

SURVEY OF CARNIVORE FAMILIES

Dogs and allies (family Canidae). *Natural history.* Canids are basically meat eaters, although some vegetable matter is taken. The gray or timber wolf (*Canis lupus*), the African hunting dog (*Lycaon pictus*), and the dhole, or wild dog of India (*Cuon alpinus*), are strictly carnivorous. The various foxes and jackals, coyotes, and the raccoon dog eat whatever food is abundant—small mammals, birds, insects, crustaceans, mollusks, fruits, or berries. The canids that are strictly carnivorous tend to hunt in packs; those that are omnivorous tend to be solitary in their hunting habits. Since the carnivorous forms depend primarily on one or a few prey species, they usually follow the moving herds of caribou or antelope, on which they feed, or move into areas where the prey species are more numerous. African hunting dogs are extremely social, always hunting in packs so intricately organized that some researchers doubt that an individual dog could survive alone. The varied diet of the omnivore reduces the necessity for organized attack on the prey species and for extensive movement.

Canids have the ability to endure a continuous chase pattern with exceptional stamina but are not capable of great bursts of speed. Their senses of smell and hearing are highly developed. Smell is used largely to track prey, and hearing to warn of impending danger. Smell is also used in association with the demarcation of territory, which tends to be outlined by frequent urination at established scent posts. Sight is less well developed, but movement is noted quickly. Canids are highly intelligent, can be trained easily, and form a basis for such sayings of man as "sly as a fox." From a behavioural standpoint wild canids are cunning or crafty, even vicious or treacherous by human standards, but usually cowardly or furtive unless running in a pack.

Canids give birth in a den in the ground, in a hollow log or tree, in a hidden brushy area, among boulders, or in a crevice of rock. The African hunting dog often dens in old diggings of the aardvark. Usually one litter of from four to six young is produced each year after a gestation period of from 51 to 80 days, depending on the species. The Arctic fox, *Alopex lagopus*, may produce two litters in a season. Most breeding takes place late in the winter, the young being born in middle or late spring. Canids living in the northern regions breed a month or so later than those in the southern regions. Carnivorous animals tend to produce young soon after the peak in the prey reproduction, a time of abundant food availability. Thus the timing of the breeding season and the duration of the gestation period must be coordinated with the season of available food in order that the young be produced when the most food is available for their rapid growth. The eyes of the young usually open in about two weeks, and they nurse for from four to six weeks. Canids of the smaller species can begin production of their own young when only one year old, but the larger forms, such as the wolf, reach sexual maturity at two or three years of age.

Canids are basically adapted to running and do so on their toes (a mode of walking or running called digitigrade). They live in a variety of habitats but generally tend to be animals of the open or grassland areas where their prey species are more abundant. Only the rare bush dog (*Speothos venaticus*) confines itself to forested areas. The red fox (*Vulpes vulpes*) tends to be an animal of the forest edge, and the gray fox (*Urocyon cinereoargenteus*) an animal inhabiting wooded areas. Thus, in North America, where both forms exist, these foxes live in slightly dif-

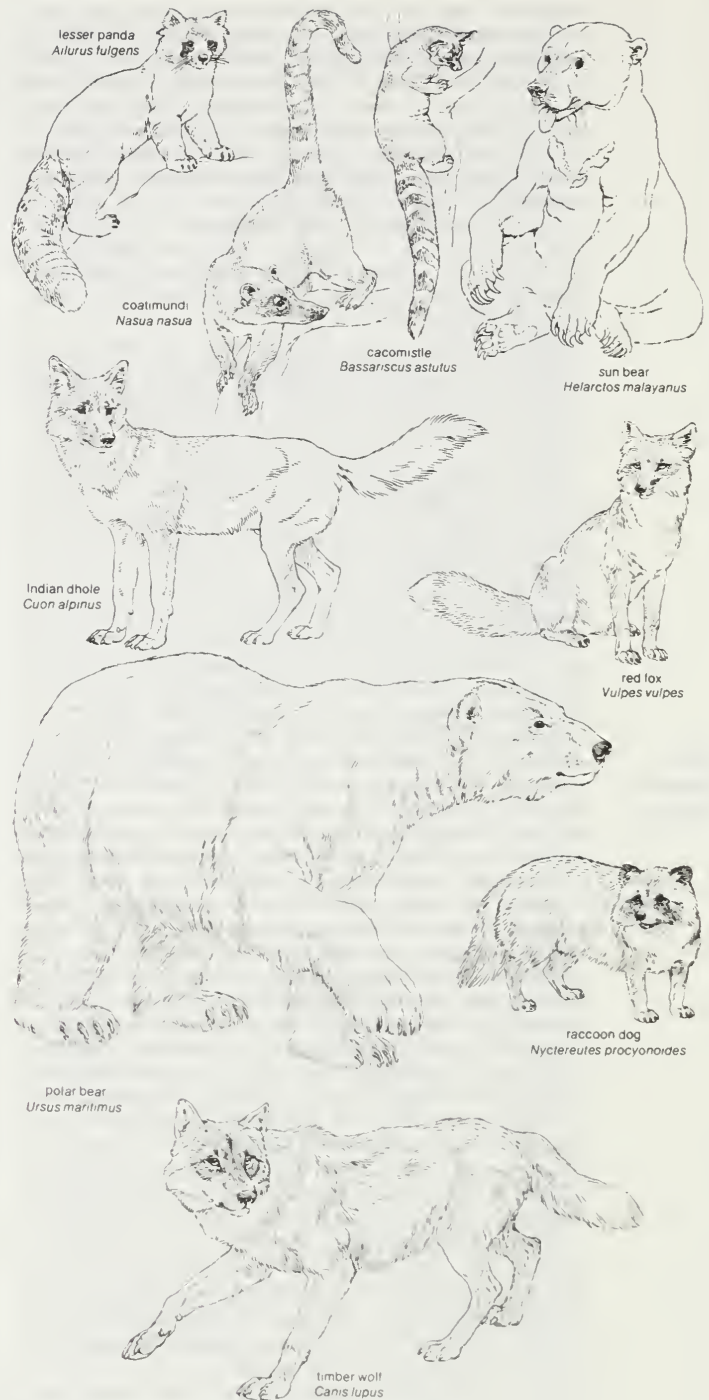


Figure 33: Representative carnivores: Procyonidae, Ursidae, Canidae.

Drawing by R. Keane. Indian dhole based on Zoological Society of London photograph.

ferent niches. Although basically terrestrial in habits, the gray fox is not averse to climbing trees. The raccoon dog (*Nyctereutes procyonoides*) often lives in tree hollows if the entrance is close to the ground. The family Canidae is worldwide, being absent only from New Zealand, Antarctica, and most oceanic islands. The Arctic fox is found farther north than any other strictly terrestrial mammal. Every major habitat has some type of canid, from the Arctic tundra, through the desert and grassland, to the tropical forest.

Canids communicate with a variety of sounds. The vocal repertoire, most highly developed in social species, includes howls, yelps, snarls, barks, and growls. These sounds are frequently associated with specialized visual signals involving movements of the ears and tail, raising of certain areas of the fur, and baring of the teeth. Within the social group or pack there is a complex dominance hierarchy,

Canid communication

Reproductive behaviour of canids

involving age levels, pair bonds, physical condition, and sexual state. Vocal and visual signals serve to minimize aggressive interactions, such as quarrels over food, that might prove injurious to the group as a whole and to the individual members. In solitary species vocalizations serve to advertise the territory, to ward off aggressors, and to communicate with the mate and young.

Importance to man. The domestic dog (*Canis familiaris*) is undoubtedly the canid of the greatest importance economically. Numerous humans are actually employed in raising, selling, doctoring, grooming, training, or providing and manufacturing food for "man's best friend." Dogs probably were the first wild animals to be domesticated; they have been found associated with Neolithic sites dating back some 8,000 years.

Some varieties of canids are important to the fur trade. At one time, a mutant form of the red fox, called the "silver fox," formed a significant part of the fur-farming industry. The raccoon dog is still raised in some regions of Russia. The more important canids whose fur is utilized are the raccoon dog, red fox, gray fox, Arctic fox, and timber wolf. Coyotes, kit and swift foxes (*Vulpes macrotis* and *V. velox*), jackals (*Canis aureus*, *C. adustus*, *C. mesomelas*), the Corsac fox (*Vulpes corsac*), the hoary fox (*Vulpes cana*), the dhole, the South African silver fox (*Vulpes chana*), the Cape fox (*Octocyon megalotis*), the maned wolf (*Chrysocyon brachyurus*), the crab-eating fox (*Cerdocyon thous*), and South American foxes (*Dusicyon* species) are occasionally used by the fur trade. Other canids, especially the red fox, are hunted primarily for sport. Some hunters call or lure foxes by using an imitation of the sound of an injured rabbit to attract the attention of the fox.

Canids also are basically important as a natural control over rodent populations. Livestock is occasionally endangered when not properly protected from the forays of wild canids. Individual foxes learn where an easy meal is located—for example, where a farmer's chickens are not properly housed. In many cases all individuals of the species are condemned because of the damaging activity of a few. Proper selective control of those individuals actually doing the damage is the only sane conservation practice. Some species also form important natural reservoirs for disease-causing organisms, especially the rabies virus. Not many canids are utilized by man for food, but the American Indian once regarded dog meat as a delicacy to be eaten on special occasions. Dog meat is also consumed in eastern Asia.

Form and function. Canids have unspecialized incisor teeth and large fanglike canines used to kill prey by slashing. The premolars are narrow and sharp and the carnassials well developed; the molars form broad surfaces that can crush substantial bones. Most canids have the basic dental formula of terrestrial carnivores, with a total of 42 teeth.

The possession of a long face or muzzle is characteristic of wild canids. All have a relatively long and obviously bushy tail. The ears are pointed, held erect, and often quite large in desert species. As well as functioning for sound detection, these large ears are believed to act as heat regulators, allowing a greater amount of heat to be dissipated in hot climates. Arctic foxes tend to have small ears, providing less loss of heat in a region where heat conservation is important to survival. Most canids have relatively long legs (especially long in the maned wolf, *Chrysocyon brachyurus*, of South America). The long legs and digitigrade gait are special adaptations for ease of running. Canids have four well-developed toes plus a dew claw (reduced hind toe) on the front foot (except in the African hunting dog, which lacks the dew claw) and four toes on the rear foot. Each toe is capped by a blunt, non-retractile claw (*i.e.*, with no sheath into which it can be withdrawn). The short claws are not adapted for use as a weapon but do aid the animal during the digging of a den or during pursuit of small prey animals in their burrows. Scent glands are often present at the base of the tail. Most canids have a uniform coloration, although there are some contrasting colours on the gray fox, a raccoonlike mask on the raccoon dog, a blotching of black, yellow, and white

colours on the African hunting dog, and a lighter-coloured belly in most canids.

Bears (family Ursidae). *Natural history.* Members of the Ursidae, except for the carnivorous polar bear (*Ursus maritimus*), are omnivorous, consuming many items that seem small for an animal the size of a bear. Ants, bees, seeds of trees, roots of the skunk cabbage, nuts, berries, insect larvae such as grubs, and even the dainty dog-toothed violet are eaten. Many bears relish honey, and the sun bear (*Helarctos malayanus*) is sometimes called the "honey bear" because of its preference for this food. Meat items taken by bears include rodents, fish, deer, pigs, and lambs. Grizzlies and Alaskan brown bears (North American subspecies of the widespread *Ursus arctos*) are known for their skillful fishing abilities during the spawning runs of salmon. The polar bear feeds almost exclusively on seals and is the most carnivorous of all bears, but little vegetation grows within its range, so the carnivorous habit is dictated by the Arctic environment. The sloth bear (*Melursus ursinus*) delights especially in raiding and destroying termite nests, sucking up termites and larvae with its funnellike lips. Bamboo shoots form the major food of the giant panda (*Ailuropoda melanoleuca*), which has a special bone formation of the forefoot that functions as a sixth digit, opposable to the other five and thus useful in handling bamboo.

Moving slowly about, bears seem to be sluggish and oblivious to much that is going on around them, but once they are disturbed their reactions may be violent and highly unpredictable. Although clumsy in appearance, a grizzly can move surprisingly fast, even through dense cover that would seriously impede a man or a horse. Their senses of sight and hearing are poorly developed, and most hunting is done by the sense of smell. Thus the smell of bacon kept in camp is an invitation to a bear in search of food. Bears, like many of their canid relatives, are solitary, except during the mating season.

Bears tend to congregate during the mating season and then pair off and mate in seclusion. Males play no role in raising the young, leaving the female soon after mating. The gestation period may vary, the fertilized egg remaining dormant in the uterus (delayed implantation), insuring the birth of young when food is abundant. Ursids breed only once a year. With a breeding season usually in late spring or early summer and delayed implantation, most young are born in January or February when the female is in the winter den. Newborn grizzlies weigh about one pound and are about nine inches long from the tip of the nose to the tip of the short tail. The Himalayan brown bear (a race of *U. arctos*) have their young later (in April and May) than the North American populations, but the polar bear has the greatest variation of any bear population in the dates when young are born (January to April). Twins are most common in bears, but up to five young may be produced. The cubs nurse for about two months and stay with the female until the next breeding (about a year and a half after birth). Most young, however, can get along on their own when about six months old. Bears reach breeding condition at three and one-half to four years of age, males usually a few months later than females. Litters are produced every second year after the initial litter.

Most bears eat large amounts of food before entering a den for a period of deep sleep during the winter. The polar bear digs a den in the snow. Grizzlies build large mounds of dirt in front of their dens. Bears are not true hibernators, since they lack the physiological characteristics (lower heart rate, body temperatures, breathing rate, and blood pressure) exhibited by those animals that do hibernate. A true hibernator is slow to be awakened from its deep sleep, and those who believe bears hibernate need only kick a sleeping bear to become convinced to the contrary. Bears usually are quiet, but they do growl at times when feeding or when being challenged by another bear.

Ursids do not roam over great distances even though they are rather large carnivores. Black bears (*Ursus americanus*) tend to stay within an area of about 36 square kilometres (14 square miles), grizzlies within an area 16 kilometres (10 miles) in diameter, or about 200 square kilometres (80 square miles). A grizzly marks the boundary of its

Distin-
guishing
physical
features of
canids

Winter
sleep of
bears

territory each spring by rubbing trees, scratching bark, or even biting large pieces from the trunks of trees. All bears, except the grizzly, Eurasian brown, and probably the polar bear, climb trees. Bears inhabit forested areas and areas that have a minimum of disturbance by man. At one time the grizzly was an animal of the plains and was one of the mammals reported by Lewis and Clark in their journey through eastern Montana in 1804. Polar bears have been noted in the open sea 65 kilometres (40 miles) from the nearest shore and on ice 320 kilometres (200 miles) from the nearest land.

Ursids are largely animals of north temperate regions, being found farther north than any other mammal (only the Arctic fox is found as far north on land). Most of the southern continents, Central and South America, Africa, and Australia, lack bears. The bear that formerly occurred in the Atlas Mountains of Morocco (*Ursus crowtheri*) became extinct in the early 19th century. The spectacled bear (*Tremarctos ornatus*) of the Andes is the southernmost form.

Importance to man. If taken when young, the grizzly can be tamed quite easily and is not uncommon in circus acts. Unfortunately, man has come to picture the bear as being quite tame and harmless, with the result that the bear is viewed without the respect that this potentially dangerous creature deserves. The grizzly is the most dangerous of the bears to man, of the most interest economically because of the damage that it can do to livestock, and of the most interest to the big-game hunter as a prized trophy. Some ursids, such as the Asiatic black bear (*Selenarctos thibetanus*), destroy fruits. At times the North American black bear does considerable damage to fruit crops and to corn fields when the corn is in the milk stage, pulling the stalks into a pile from as far as the bear can reach.

The pelts of bears have been used for a number of purposes. Perhaps most popular has been the bearskin rug. Skins also have been used for lap rugs or sledge rugs, hats for the English Guard regiments, trimmings of coats, and for muffs. Skins of the rare giant panda commanded a high price in the early 1900s, but this animal apparently is close to extinction today. Black bears are prized as a food source in China, as is the flesh of the polar bear in the far north. The liver of the polar bear, however, is highly poisonous, due to high levels of vitamin A. The teeth and claws of bears have been favourite ornaments with the Indians and Eskimos for years, and the fat furnishes the "bear grease" of commerce.

Form and function. The teeth of the omnivorous bears are unspecialized. The first three premolars are usually either missing or extremely small. The carnassials are poorly developed, and the molars have broad, flat crowns. Except for variability as to the presence of premolars, the ursid dental formula is that of the Carnivora generally. The sloth bear lacks one pair of upper incisors. The jaw of the bear is controlled at the hinge by a powerful set of muscles.

Bears have an elongate skull that is especially heavy in the back portion. These plantigrade animals (plantigrades walk on the sole and heel of the foot) are powerful in build and have relatively short, massive legs with five toes on each foot. The toes end in enlarged, nonretractile claws that are especially well developed in the sloth bear and to some extent in the grizzly. The claws on the front feet are usually better developed than those on the rear and are especially adapted for digging out marmots or other small rodents. Bears, however, seldom use the claws to construct burrows. They generally have naked soles, but those of the polar bear are covered with hair, enabling the animal to walk on ice with a firm footing. In swimming, the polar bear uses only its front limbs, an aquatic adaptation found in no other four-legged mammal. Bears lack a clavicle but have a baculum (penis bone). Their lips, which are not attached to the gums, are protrusible and mobile. Ursids have a short (7 to 13 centimetres; 2.8 to 5.1 inches), stubby tail and are the largest living terrestrial members of the Carnivora. The Alaskan grizzly may weigh as much as 780 kilograms (1,700 pounds), and the smallest bear, the sun bear, weighs approximately 27 kilograms (60 pounds). In most species the two sexes are about the same size, but the male tends to be larger in

the grizzly and polar bear. Bears of the genus *Ursus* tend to be of a uniform colour, ranging from white to brown, blue gray, and black. Unlike the majority of carnivores, ursids are often black or mostly black. The North American black bear occurs in a number of colour variants or phases; in addition to the typical black phase, there is a tan-brown phase (called the cinnamon bear) and a bluish-black phase (called the glacier bear). The spectacled bear and the several specialized Asian bears are black with white or tan markings. The giant panda is strikingly black and white, and the remaining genera largely black, with areas of white or orange on the chest, neck, or face. The fur is coarse, long, and shaggy in most northern forms.

Raccoons and allies (family Procyonidae). *Natural history.* The family Procyonidae includes raccoons, coatis or coatimundis, olingos, and the ringtailed cat, kinkajou, and lesser panda. Procyonids, like bears, are omnivorous, eating numerous types of insects, crayfish, crabs, fishes, amphibians, reptiles, birds, small mammals, and a variety of fruits, nuts, roots, and young plants. A characteristic of most omnivorous animals is that they take whatever foods are available, varying their food habits with the season, the locality, and the abundance or population levels of the food items from one year to the next. The ringtailed cat (*Bassariscus astutus*) tends to be more carnivorous than the raccoon (*Procyon lotor*). The lesser panda (*Ailuus fulgens*) is the most vegetarian of the group.

The various senses of procyonids do not seem to be especially well developed, but the raccoons and lesser panda appear to be above average in intelligence and are often described as being cunning. Raccoons are well known for their insatiable curiosity and ability to use their "hands" in performing feats such as opening doors and getting into mischief when kept as pets. This curiosity of raccoons is well known by trappers, who have used bright objects placed on the treadle of a trap to attract them. The sense of touch probably is of the greatest importance to the raccoon in its food-gathering activities. Most procyonids, except perhaps the coatis (*Nasua* species), are more active in the twilight and evening hours than in the daylight. Procyonids tend to be solitary, although families often move and feed as a group. Coatis travel in bands, some individuals moving through the trees as others follow on the ground. Olingos (*Bassaricyon* species) also move about in small bands. All procyonids are arboreal and move around in trees with agility. On the ground, the raccoon resembles a small bear as it ambles along.

Procyonids breed in late winter or early spring. During copulation, the female raccoon makes sharp rattling cries. After a gestation period of about two months (for the raccoons 54 to 66 days, for lesser panda 90 days), two to six young are produced. One litter a year is usual. Young raccoons "twitter" like young birds when disturbed in the nest, and the young of other procyonids also produce high-pitched calls. Most adult procyonids produce a variety of snarls, growls, whines, screams, and barks. Most of these calls have little carrying power. Young raccoons are dependent on the female for the first 10 weeks but may nurse until approximately 20 weeks of age (ringtails to about four months) and may stay in the family group into winter. Being mainly arboreal, young procyonids are born in tree dens. Raccoons use old buildings or dens in the ground if tree dens are scarce. Ringtails use crevices in rocks.

Except for the lesser panda, which is found in the southeastern Himalayas of Asia at elevations from 1,980 to 3,650 metres (6,500 to 12,000 feet), procyonids are found naturally only in the New World. They are largely mammals of the lower elevations of temperate and tropical regions, in areas well supplied with water and trees. Raccoons have been introduced successfully into many parts of Europe (Germany, White Russia, the Ukraine) and Asia (Lake Balkhash, Siberia).

Importance to man. The raccoon is the only procyonid of major economic significance. Raccoon hunting has been a sport of substantial importance in the United States, and in some parts of the United States many people own dogs trained to hunt raccoons. The coatimundi also is hunted with dogs and its flesh is also eaten. Meat of the raccoon is

Use of
bear pelts

Bear claws

Procyonid
breeding
habits



Figure 34: Representative carnivores: Mustelidae.

Drawing by R. Keane, striped skunk based on E.P. Walker, *Mammals of the World*

in more demand as food than that of any other carnivore. The meat is roasted, fried, stewed, barbecued, fricasseed, and made into patties. The "coonskin" cap has become part of American legend. The fur of the raccoon is one of the staple items of the fur trade in the United States today, as it has been over many of the past 100 years. Also used are the furs of the ringtailed cat, kinkajou (*Potos flavus*), and the lesser panda.

Raccoons at times do damage to fields of corn and to abandoned buildings, which they frequently inhabit when natural dens are scarce. In certain regions, raccoons become a problem to nesting waterfowl, especially wood ducks and goldeneye, which, like the raccoons, raise their young in hollow trees.

Most procyonids make interesting pets. Raccoons, kinkajous, and coatis are frequently sold as house pets, but their inquisitive nature and climbing ability make them troublesome if allowed to run loose indoors.

Form and function. The unspecialized teeth of procyonids are adapted for an omnivorous diet. The elongate canines are oblong in cross section, and the premolars are small, narrow, and pointed. As in bears, the carnassials are poorly developed and the molars are broad. Procyonids have the basic dental formula of the Carnivora, except that there is one less lower molar.

Procyonids are plantigrade or semiplantigrade, with legs of moderate length and five flexible toes with nonretractile or (in the ringtail cat and lesser panda) semiretractile claws. The tail is well developed (20 to 70 centimetres [7.9 to 27.6 inches] long), bushy, and in all forms except the kinkajou ringed by light and dark bands. In the kinkajou, it is prehensile, gripping a branch like a fifth limb. The anal scent glands of the kinkajou and lesser panda produce a musky odour when the animal becomes excited. The snout is extremely flexible in the coati and somewhat movable in other procyonids. The members of this family are medium in size, ranging from 0.8 to 22 kilograms (1.75 to 48.4 pounds) in weight (the record weight of a raccoon is 30.2 kilograms [66 pounds, 6 ounces]). There is no difference in size between males and females. In coloration the procyonids range from gray or brown to the rich, reddish brown of the lesser panda.

Mustelids (family Mustelidae). *Natural history.* The family Mustelidae contains a variety of animals unmatched by any other family in the Carnivora except the civets (Viverridae). The family includes the weasels, ferrets, mink, marten, fisher, skunks, wolverine, otters, badgers, and a number of less well-known animals, a total of about 70 species in 25 genera. The weasels, mink, and polecats (*Mustela* species) are fierce predators, feeding exclusively on small mammals and birds, hunting and trailing their prey with a keen sense of smell. At times, bird eggs, salamanders, snakes, and insects are taken. The smallest mustelid, the least weasel (*M. nivalis*), consumes a third to half its body weight per day, killing the prey by biting at the base of the

skull. Weasels often store food in underground caches. The Asiatic polecat (*M. eversmanni*) and the North American black-footed ferret (*M. nigripes*) feed almost exclusively on the animals whose colonies they inhabit, the polecat on marmots and ground squirrels, the ferret on prairie dogs. Both of these animals occasionally feed on insects. Mink (*M. vison* and *M. lutreola*) feed largely on aquatic animals (muskrats, crayfish, mollusks, and fish) as well as on organisms inhabiting the banks of bodies of water.

The fisher (*Martes pennanti*) feeds mainly on mammals, but a few nuts, birds and eggs, fish, and seeds are taken. This animal is noted for its ability to kill and eat porcupines. Small terrestrial mammals make up over 90 percent of the food of the American marten (*M. americana*). The Eurasian pine marten (*M. martes*) and the stone or beech marten (*M. foina*) feed largely on small mammals but in the autumn may turn to wild berries and fruit. The yellow-throated marten (*M. flavigula*) of Asia is one of the largest martens and feeds on the young of larger mammals such as musk deer, wild boars, roe deer, spotted deer, and raccoon dogs, as well as on squirrels, rabbits, and the smaller rodents. These large mustelids also take a variety of other food items: mollusks, grasshoppers, spawning fish, birds and their eggs, cedar nuts, grapes, and berries.

Badgers are quite variable in their diet. About 50 percent of the diet of the American badger (*Taxidea taxus*) is composed of ground squirrels, which are captured in their burrows. Badgers can dig with remarkable speed and can overtake the rodent when the latter has taken refuge in a blind tunnel. The remainder of their food is made up of small rodents, rabbits, birds and their eggs, snakes, turtle eggs, and insects. Insects, in fact, form a staple food of the young badger. The Eurasian badger (*Meles meles*) and the ferret badgers (*Melogale* species) are more omnivorous than the American badger and prefer insects (especially bees) and their larvae (grubs). Earthworms, mollusks, crustaceans, frogs, lizards, birds, eggs, as well as vegetable matter, such as berries, nuts, grains, fruits, bulbs, and green vegetation, are also taken. Rats or honey badgers (*Mellivora capensis*) eat small mammals, lizards and snakes, insects, fruit, and especially honey. An interesting association exists between the African honey guide (a piciform bird) and the honey badger. When the bird finds honey it gives a special call, which attracts the badger. The latter breaks open the hive with its long, sharp claws, eats the honey, and leaves the wax and bee larvae for the bird. If the bird does not succeed initially in attracting a honey badger, it seeks one out and, using the "honey" call, which the badger recognizes, leads it to the beehive.

Skunks (*Mephitis*, *Spilogale*, *Conepatus*) are omnivorous mustelids, with small mammals forming an important part of the diet. Skunks may be becoming scarce in some parts of North America due to poisons gradually accumulated through extensive feeding on insects poisoned in agricultural control programs.

Otters feed largely on the invertebrates and vertebrates (especially fish) found in their aquatic environment. The oriental small-clawed otter (*Aonyx cinerea*) and the clawless otter (*Aonyx capensis*) feed to a greater extent on invertebrates (clams, snails, crabs) than does the river otter (*Lutra* species). *Lutra canadensis* at times takes large numbers of crayfish. The African small-clawed otter (*Aonyx* species) is believed to feed on small mammals, the eggs of birds, and amphibians more than on fish. The sea otter (*Enhydra lutris*) is perhaps known best for its habit of breaking open sea urchins and clams on a rock held on its chest while it floats on its back. Seaweed, cuttlefish, and small fishes also are eaten.

Sharp
senses of
mustelids

Smell apparently is the most important sense to mustelids hunting on land. Black-footed ferrets have been noted to stop repeatedly and sniff the air while hunting. Hearing is also well developed and seems to be utilized as a means of detecting danger. Sight is less developed and, in the ferret at least, of little importance beyond about 90 metres (300 feet). Of most importance to otters, living in an aquatic environment, is an acute sense of touch, through their very sensitive whiskers, which they use as an aid to guide them through the water. Most mustelids are silent, seldom making more than a squeak, whistle, or bark. Many will snarl or growl when annoyed. The giant otter (*Pteronura brasiliensis*) of South America is perhaps the most vocal, with its high-pitched whistle. Yellow-throated martens give a harsh cry when excited and often chuckle in a low tone.

Almost all mustelids are active both day and night, although most of their activity is nocturnal. Only the Old World badgers (*Meles meles*), ferret badgers (*Melogale* species), spotted skunks (*Spilogale* species), and hog-nosed skunks (*Conepatus* species) seem to be almost strictly nocturnal. The giant otter and sea otter seem to be more diurnal than other mustelids.

Except when travelling with young in a family group, most adult mustelids are solitary in habits. Sea otters seem to be more or less gregarious, as are Old World badgers. Badgers build a labyrinth of underground passages that some call "badger cities" because of the extent of the burrows and number of badgers that inhabit them. Female skunks (*Mephitis* species) often den together in large groups during the winter. Grisons (*Galictis* species) and yellow-throated martens often travel in groups of five or six. Tayras (*Eira barbara*) and ratels tend to travel in pairs.

Since a great variation exists in the habits of mustelids, a great diversity would be expected in the habitats in which they live. Most mustelids, however, are terrestrial, living in forested or brushy areas. Many species can survive in a variety of habitats, from forest to desert. Some, like the marten and fisher, are strongly arboreal. Others, like the otter, are largely aquatic. The New World skunks and the North African spotted weasel (*Poeciliotis libyca*) tend to inhabit areas farmed by man. The ferrets, American badger, and Eurasian polecats live mainly in grassland or semi-arid areas. Wolverines live in the more northern areas of taiga and forested tundras.

The locomotion of mustelids also varies. Most small mustelids scamper across the ground, making intermittent bounds as they move along. The larger forms sometimes lumber slowly. During their travels, mustelids tend to stop, sit up on their haunches, and look over the area. Some are extremely agile in trees; others are excellent swimmers.

Delayed implantation is perhaps more characteristic of the mustelids than any other group of the Carnivora. Members of the genus *Mustela* have gestation periods (determined from the date of fertilization) ranging from as much as 220–337 days (long-tailed weasel, *M. longicauda*) to as little as 36–42 days (ferrets). In the genus *Martes* the period is highly variable, between 220 and 297 days. The gestation period of badgers (*Taxidea*, *Mellivora*, *Meles*) ranges from six to eight months and that of otter from eight months to a little more than a year. Skunks (*Mephitis* species) do not have delayed implantation and apparently have the least variation in the period of gestation (62 to 72 days). Mink (*Mustela vison*) usually produce young in 45 to 50 days, with extremes from 38 to 76 days. On mink ranches, the later in the year that mink are bred, the shorter the gestation period.

The breeding season usually is in the late winter or early spring (mink, skunks, polecats, otters) or in the summer (weasels, marten, wolverine, sables, European and American badgers). Some mustelids breed at almost any time of the year (least weasel, river and sea otters). Those with exceptionally long gestation periods breed soon after a litter is born or in the following year. Some forms have two litters a year (striped skunk, least weasel, ratel), but in most only one litter is produced. Average litter size varies from one young in the sea otter to between eight and 12 young (up to 18) in the Asiatic polecat. Most mustelids have between three and five young in a litter. The young of mustelids are born in ground dens, in rock crevices, under the roots of trees, under haystacks or buildings, or in hollow trees (the pine and American martens). The young of most species are weaned when between six and 10 weeks of age, but young river otters nurse for about four months. Wolverine young are allowed to remain with the female for two years before establishing their own territories. In most species, however, the family disperses near the end of the summer or early fall. Mustelids are found worldwide, except for most oceanic islands, Australia, and Madagascar. Introductions have been made on some of the larger islands, such as New Zealand. Mustelids, like bears, generally are more abundant in the more northern continents.

Importance to man. Of all mammalian families, the Mustelidae contains the greatest variety of valuable mammals utilized by the fur trade. Fur coats that command the highest prices are usually those made from these luxuriant furbearers, especially mink, ermine (short-tailed weasel), and sable (a name applied by the fur trade to several species of *Martes*, especially the Eurasian *M. zibellina* but also the American *M. americana*). Excellent skins are produced by the marten, fisher, skunk, and otter. In fact, the durability of all furs, including those from other mammalian orders, is based on the river otter, which has been given the top rank of 100. Pelts of the ermine and sea otter have been the furs of royalty for hundreds of years. More than \$100,000,000 worth of ranch mink skins alone are produced yearly, and some breeds (like the Kojah mink) brought \$2,700 per skin in the late 1960s. During the 18th and 19th centuries a collar made from the sea otter was regarded by some as the highest of status symbols, equivalent to owning a Rolls Royce or private yacht today. In 1968, at the first sale of sea-otter skins in 55 years, the better-quality skins brought \$2,300 apiece.

The pelt of the wolverine is prized highly by the Eskimos and others in Arctic areas because, when used for the lining of a parka, it can be cleared of frost from breath by a sweep of the hand. Moisture does not freeze to wolverine hair as it does to other furs. Badger hair has been utilized in the past to manufacture shaving brushes. Hair from the tails of many species is in high demand for paint brushes. Among the Chinese, for example, the tail of the kolinsky or China mink (*Mustela sibirica*) is used for making delicate paint and artists' brushes and is in more demand than is the rest of the pelt. The tails are removed from the skins and sold by the pound. Hairs from the tail of the sable also make excellent brushes, because the taper of the individual hairs causes the brush to form a sharp point when wet.

Being largely carnivorous, the mustelids form an important link in the biological food chain. Rodents are kept at a reasonable level by this natural control. The spotted skunk, ferret badger, long-tailed weasel, Patagonian weasel, zorille (*Ictonyx striatus*), and sometimes the striped skunk (*Mephitis mephitis*) have been encouraged to live in or near the habitations of man because of their ability to keep the rodent and pest insect populations under control. Systematic destruction of these valuable predators can be a costly mistake to farmers.

Being carnivorous, mustelids sometimes become pests, especially around poultry—particularly when the poultry farmer is lax in his sanitation. Weasels, marbled polecats, and striped skunks are the common offenders. Considerable exaggeration of the prowess and deliberate maliciousness of the wolverine exists in literature. The behaviour of this animal has led many trappers to regard it as a

Mustelid
reproduc-
tion

Importance
in the food
chain

true "demon of the north." A wolverine travels a trapline because it learns that a ready source of food is being held captive. In many cases the wolverine will follow the trapline back to the trapper's cabin and ransack the area in its search for food. The largest of the mustelids (around 11.5 kilograms, or 25 pounds), this animal can do extensive damage.

Those mustelids that are proficient diggers (e.g., badgers) sometimes become a hazard to horses (and riders), which stumble when stepping into a burrow. Other mustelids, such as the North American skunk, tend to be one of the more common wild mammals that serve as a reservoir for rabies. Meat of the Malayan stink badger (*Mydaus javanensis*), and sometimes the skunk (*Mephitis mephitis*), is eaten in some areas of the world. The obnoxious scent glands of mustelids discourage man from eating most of these animals, however. On the other hand, the scent from these glands was used in the past as a base for perfumes. Grisons and the ferret (*Mustela putorius*) have been trained by man to enter the burrows of other animals and drive them out; the grison to flush chinchillas, and the ferret, rabbits and rats. In China Oriental clawed otter and river otter are trained to catch fish or to chase them under nets which are then dropped, securing the catch. In some areas the bones or various parts of the tayra, African striped weasel, Malayan stink badger, hog-nosed skunks, and the African small-clawed otter are believed to possess medicinal value.

Form and function. Being highly carnivorous, mustelids have well-developed carnassials and a reduced number of premolar and molar teeth, the total number of teeth being between 30 and 38 (average, 34). With such a reduction in tooth number, the facial area of these animals is shortened. In the wolverine and American badger the bone of the lower jaw has a projection, the postglenoid process, that curves over the jaw articulation (glenoid fossa), resulting in a condition in which the mandible can be locked in place.

Mustelids are either digitigrade or plantigrade, and they have five toes on each foot. Most mustelids have short legs and many tend to have an elongate, slender body. They are quiet, agile, and graceful in their movements. There is a great variation in the tails of mustelids: flattened and elongate in the otter; stubby in the stink badgers; short in some of the weasels and badgers, long and thin in other weasels; long and bushy in marten; and especially elongate, with long hairs, in the skunks. The ears are small and rounded. Colours are usually brown or black. Some weasels change colour, becoming white in the winter, thus protecting the animal or, perhaps more probably, hiding it from its prey and making capture easier. Such white weasels are called ermine in the fur trade. Most mustelids have well-developed anal scent glands. The skunks, zorrillo, and the marbled polecat (*Vormela peregusna*), and to some extent the stink badger and ratel, can forcibly eject the vile-smelling material as a spray or fluid. All mustelids that possess powerful, controllable scent glands (except the stink badger) are marked in obvious contrasting colours, usually black and white. Such a striking colour combination may serve as a warning to potential predators, minimizing the likelihood that the animal will be attacked. After a single encounter with an adult mustelid, a predator avoids all animals so marked, including the less formidable young of the mustelid.

The smallest member of the Carnivora, the least weasel, belongs to the family Mustelidae and weighs 30 to 70 grams (one to 2.5 ounces). The largest mustelid is the sea otter, which may weigh 41 kilograms (90 pounds). In most mustelids the males are larger than the females.

Civets and allies (family Viverridae). *Natural history.* Viverrids are more numerous than the mustelids in terms of living genera (about 36) and species (about 75).

They are largely omnivorous in food habits, eating a variety of small mammals, birds, reptiles, eggs of birds and reptiles, amphibians, fishes, crustaceans, insects and their grubs, and earthworms, as well as vegetation such as fruits, nuts, bulbs or roots, and other plant material. A few, such as the palm civets (*Nandinia*, *Arctogalidia*, *Paradoxurus*, *Macrogalidia*), tend to be mostly vegetar-

ian but at times take small mammals, birds, and insects. The binturong (*Arctictis binturong*), the largest member of the family, feeds mainly on fruit, although carrion is sometimes eaten. Although the suricates or slender-tailed meerkats (*Suricata suricatta*) tend toward an omnivorous diet, bulbous roots constitute most of their food. Approximately half of the genera of viverrids are composed of species that tend to be more carnivorous than herbivorous; the remaining genera are truly omnivorous. Some species tend to feed on one type of food: the ruddy mongoose (*Herpestes smithi*) on large snails, the Indian mongoose (*H. edwardsi*) on carrion remaining from the feasts of the large predators, the water mongoose (*Atilax paludinosus*) on crocodile eggs, the striped-necked mongoose (*H. vitticollis*) and the Congo water civet (*Osbornictis piscivora*) on fish. The banded palm civet (*Hemigalus derbyanus*) and Owston's civet are sometimes known to eat significant amounts of earthworms.

Drawing by R. Keane

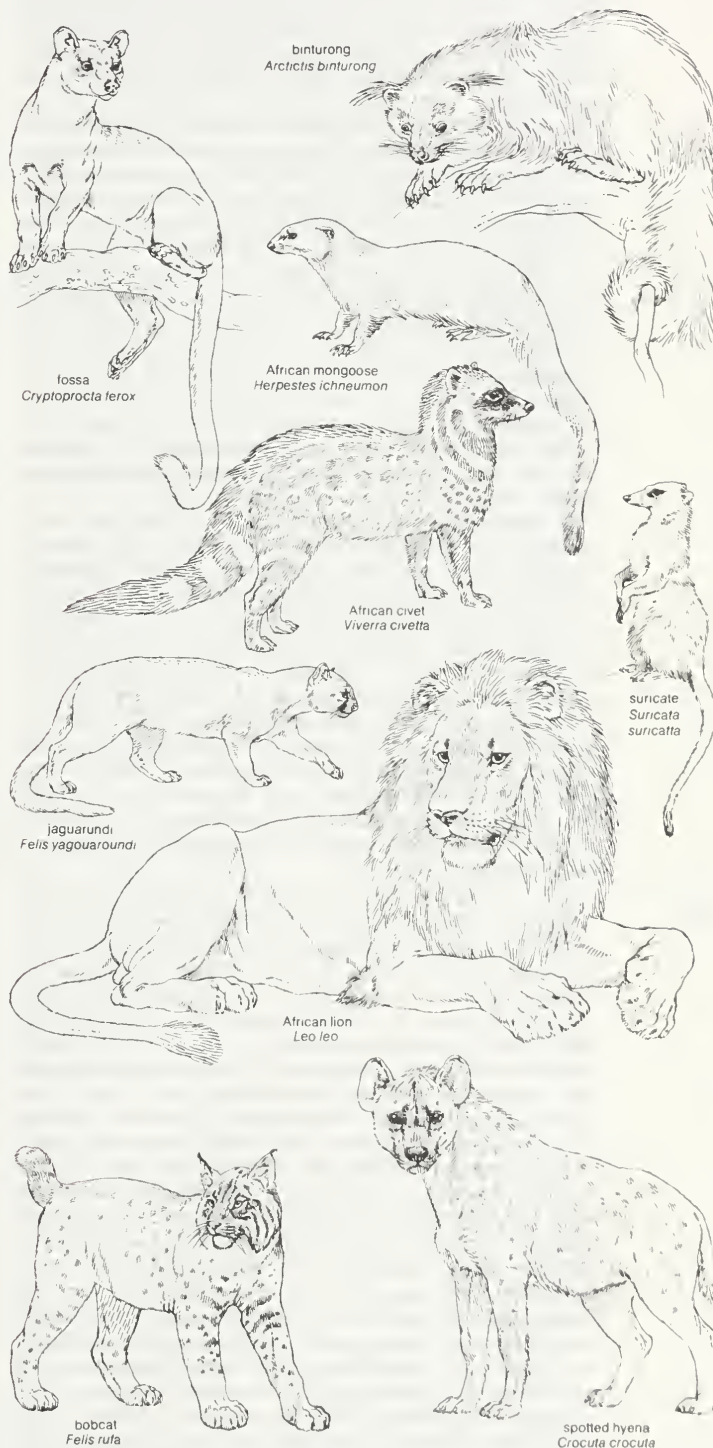


Figure 35: Representative carnivores: Viverridae, Felidae, Hyaenidae.

Repellent
scent of
mustelids

A number of viverrids (*Herpestes*, *Atilax*, *Mungos*, *Helo-*
gale) are noted for their peculiar habit of opening eggs, as
well as other food items with hard shells (crabs, mollusks,
nuts). The animal stands on its hindlegs and hits the egg
against the ground. Sometimes it carries the egg to a rock
and, standing with its back to the rock, throws the egg be-
tween its legs and against the rock until the shell is broken.
(Early reports of this behaviour met with skepticism but
have been verified by other observers.) The Madagascar
narrow-striped mongoose (*Mungotictis lineatus*) exhibits
the same behaviour but lies on its side and uses all four
feet to toss the egg.

Viverrids
as snake-
killers

Viverrids are perhaps best known for their ability to kill
some kinds of poisonous snakes, especially cobras, which
are slow in striking and apparently unable to strike when
at close quarters. The mongoose can dart within the outer
striking range of the snake and grab and break the lower
jaw. The mongoose is much less successful against a snake
that has a very rapid strike and return to the pre-strike
defensive position, such as a rattlesnake. The white-tailed
mongoose (*Ichnemumia albicauda*) has the unusual habit of
feeding on the tree hyrax, which it follows and kills in the
tops of trees.

Viverrids are highly variable in hunting and social be-
haviour. Some species are solitary; others hunt in groups,
ranging from a single pair to a family group or a wander-
ing band of more than 20 individuals. Two species, the
meerkats (*Suricata suricatta* and *Cynictis penicillata*), live
in colonies. Colonial organization is unusual for animals
that are largely carnivorous but is found also in the Euro-
pean badger. An interesting relationship exists between
the type of association and the time of day that viverrids
are active. The viverrids that travel singly or hunt in pairs
comprise the majority of species, and all are nocturnal.
All those viverrids that hunt in bands or live in colonies
are diurnal. The only exceptions are the mongooses of
the genus *Herpestes*, which may hunt singly or in bands
and may be either nocturnal or diurnal depending on the
species involved.

The senses for seeing, hearing, and smelling are well
developed in viverrids. Predatory birds apparently are rec-
ognized almost instantly and at considerable distances.
Viverrids, like mustelids, tend to be silent, but such sounds
as purring, grunting, chirping, low-pitched coughing,
growling, hissing, whining, twittering, and barking have
been described. More elaborate or complex "language"
usually is found in gregarious forms. Specific warning or
danger cries, signifying the approach of either a ground
or aerial predator, are known in the social meerkats. The
noisiest viverrid, however, is the nocturnal, arboreal bin-
turong, which growls, hisses, and, at times, howls loudly.
Also extremely noisy are the cusimanses (*Crossarchus*
species), which travel in large groups grunting, chattering,
and twittering.

Viverrid
habitats

Most viverrids are terrestrial, but at least 10 genera have
forms that inhabit trees to a considerable extent. Mostly
arboreal are the African palm civet (*Nandinia binotata*),
the small-toothed palm civet (*Arctogalidia trivirgata*), the
musang (*Paradoxurus* species), the masked palm civet
(*Paguma larvata*), the binturong, the fossa (*Cryptoprocta*
ferox), and the slender mongoose (*Herpestes sanguineus*).
The linsangs (*Prionodon* species), although spending some
time on the ground, are good climbers and live in the
hollows of trees. Palm civets (subfamily Paradoxurinae), a
few other civets, and genets also spend considerable time
in trees. Mongooses, however, seldom climb and are un-
able to move about easily in trees. This inability is evident
when they descend from trees. An efficient well-adapted
climber comes down a tree head first. Of the mongooses,
only the slender mongoose, and perhaps the Indian gray
mongoose, have been observed descending head first. All
other mongooses back down. Although most viverrids can
swim if need be, three genera have semi-aquatic members:
the Congo water civet, the African civet (*Viverra civeta*),
and at least three species of mongooses of the genus *Her-*
pestes (*H. brachyurus*, *H. urva*, and *H. vitticollis*). The
otter civet (*Cynogale bennetti*) has slightly webbed feet
and is aquatic. An even more aquatic viverrid, or at least
a more agile one, is the water mongoose, which dives and

swims like an otter, although, unlike the otter, its feet are
not webbed.

Some mongooses exhibit curious antics to attract their
prey. For example, the white-tailed mongoose may dance
outside chicken pens and bite off the heads of chickens
that curiously stick their heads through the netting. The
banded mongoose is said to stand upright and fall over
on its side, using such antics to attract curious guinea
fowl so that they approach close enough for capture. An
interesting habit found in the meerkat plays a role in heat
regulation. The belly of this animal has little hair, much
less than the back. In the winter or during cold weather,
the meerkat basks in the sun, standing on its hind feet
with the ventral surface directed toward the heat source.
When the temperatures are high, the animal seeks out cool
surfaces, on which it lies on its stomach. It may even dig
out an area to expose the cool subsurface layers of earth
before lying down.

The fossa walks on the soles of its feet and thus is planti-
grade or semiplantigrade, but most other viverrids walk
on their toes (digitigrade).

Information concerning reproduction in viverrids is
sketchy. In most species studied there are two litters a
year (usually one in spring and one in fall), although
breeding can occur throughout the year. The gestation
period is usually between 49 and 64 days, but in various
mongooses of the genus *Herpestes* it ranges from 32 to 49
days. Most litters consist of three or four young, but the
black-legged mongooses (*Bdeogale* species) may have only
a single young. Most female viverrids have two pairs of
abdominal nipples, thus perhaps restricting the litter size
to four, but some females of the genus *Herpestes* have
three pairs. The young are born in burrows, in a nest of
grass on the ground, or in hollow trees.

Viverrids are found in the Old World tropics and sub-
tropics, throughout Africa, southern Asia, Madagascar,
and southeastern Europe. Only the genet (*Genetta genetta*)
occurs naturally in Europe. Although restricted to south-
ern France and Spain today, it was at one time found
as far north and east as Belgium and Germany. Various
mongooses (*Herpestes* species) have been introduced into
the West Indies, Hawaii, and New Zealand, as well as Italy
and certain areas of the Balkan Peninsula. Viverrids are
most abundant in Southeast Asia (India and the Malay
Peninsula). In habits and habitat they fill the niches of
the mustelids and procyonids found more abundantly in
the northern regions or in the New World. Viverrids are
the only members of the Carnivora to reach Madagascar,
where they have flourished and undergone adaptive radi-
ation. About a sixth of the known genera of viverrids are
restricted in distribution to this large island.

Importance to man. Probably of most significance eco-
nomically is the musk, called civet, produced by the anal
gland in members of the viverrid group known as "civet
cats." Civet is a yellowish substance with the consistency
of butter and is of great importance as a base for per-
fumes. It is produced by the anal scent glands and stored
in a sac near the genital region. The Oriental and African
civets (*Viverra* species) and the lesser Oriental civet, or
rasse (*Viverricula indica*), are kept in captivity so that this
musky secretion can be collected several times weekly. The
secretion is obtained by scraping the sac with a wooden
spatula. As much as four grams of civet can be obtained
in one week. Long, narrow cages are used to prevent the
animal from turning and biting while the secretion is be-
ing collected.

This musk has been in great demand for many years
and has commanded high prices. Civet is used as a base
for perfume extracts, and in India the product of the
lesser Oriental civet has been used to flavour tobacco. In
some areas, parts from viverrids (the binturong and the
linsang, *Poiana richardsoni*) are believed to have various
medicinal values.

The only civets of value to the fur trade are certain mem-
bers of the subfamily Viverrinae. Of these, the large In-
dian, or "Chinese," civet (*Viverra zibetha*) is the only one
to be utilized for fur to any great extent. Fur of the small
Indian civet and the genet is used for trimming clothing.
Some members of tribes in Madagascar use the tails of

Viverrids
as the
source of
musk

civets for ornamentation. The meat of both the large and small Indian civets has been described as delicious.

The genet, white-tailed mongoose, and fossa often become pests in areas where poultry is raised. Most viverrids are efficient rat catchers, and some, such as the binturong, suricate, cusimanse, and certain mongooses, are domesticated because of their affectionate nature and their ability to eliminate mice, rats, and cockroaches from the residence of the owner. In some areas, Greece and southern Italy for example, genets have been kept as house pets, and this animal may have been the "cat" of the ancient Greeks. In other areas, the binturong is tamed to follow its master like a dog.

Form and function. Viverrids have from 36 to 42 teeth. The canines are slender and elongate and the carnassials well developed.

Viverrids have short legs with five toes on each foot, except for the meerkats; *Suricata* has four toes on each foot and *Cynictis* four on the hindfoot. The pollex ("thumb") and hallux ("big toe") are functionless and do not touch the ground. Each toe is capped with semiretractile claws that are more retractile in digitigrade species than in those that are semiplantigrade. Most viverrids have a long, slim body and a long, bushy tail. Except for some marsupials, the binturong is the only Old World mammal with a prehensile tail. The elongate head ends in a pointed, foxlike muzzle. Coloration in the civets can be quite variable in pattern, from uniform to spotted, blotched, or striped, and ranges from black and gray through various shades of brown to rufous and yellow. The stripes may run longitudinally (*Mungotictis*, *Galidictis*) or vertically (*Mungos*, *Chrotogale*, *Hemigalus*). Some forms of civets (especially *Bdeogale*, *Ichneumia*) have black legs.

Viverrids range in size from the dwarf mongoose (*Hologale parvula*) with a weight of about 500 grams (about one pound) and a total length of 300 millimetres (12 inches) to the binturong with a weight of up to 14,000 grams (about 30 pounds) and a length of 1,800 millimetres (72 inches).

Hyenas (family Hyaenidae). *Natural history.* Hyenas form a relatively homogeneous family that includes two basic types: the true hyenas (two species of *Hyaena* and one of *Crocota*) and the aardwolf (*Proteles cristatus*). Hyenas are noted for their scavenger feeding and bone-crushing habits, although they do hunt live animals. Their food consists largely of the remains of artiodactyls, such as the antelope, left by the large cats. Vultures and jackals often feed on such kills after the cats and before the hyenas move in. The spotted hyena (*Crocota crocuta*) eats large amounts of bone and skin, seeming to prefer them to meat. These, of course, are the remains after the predator has had its fill. Spotted hyenas hunt in aggressive packs and at times attack and kill larger animals, even as large as a young rhino. If the animal attacked puts up a fight, the hyenas stop the attack and move away. In addition to carrion, brown hyenas (*Hyaena brunnea*) feed on small mammals, young or newborn larger mammals, reptiles, eggs, insects, fruits, and berries. In coastal areas they may feed on dead crabs, whales, and other sea life washed up on shore. The brown hyena is less apt to attack living animals than is the spotted hyena. The striped hyena (*Hyaena hyaena*) has food habits similar to those of the brown hyena. The aardwolf feeds almost exclusively on termites but eats a few other insects, such as beetles. A large quantity of soil and grass is usually eaten accidentally while feeding at termite mounds.

Since carrion produces a distinctive odour, hyenas depend on the sense of smell to a great degree. Most of their food probably is located by this sense alone, sight and hearing being less important. Sight is poorly developed in hyenas.

The well-known howl of the "laughing hyena" is the call of the spotted hyena. The other species are not as noisy.

In areas undisturbed by man, hyenas hunt during the day. Members of this family inhabit open areas of plains and brushland but at times may be found in open forested areas. Hyenas form packs and are tireless trotters and runners, at times wandering great distances.

Spotted hyenas breed during the dry months, from April to July; the aardwolf breeds from October to November.

The gestation period is approximately three months (95 to 110 days in the spotted hyena). Hyenas usually have three or four young. Dens may be located in abandoned burrows, caves, dense brushy areas, or among rocks. Some hyenas use the burrows of the aardvark for nesting spots, as does the African hunting dog.

The Hyaenidae have a distribution similar to, but more restricted than, that of the viverrids. They are found throughout Africa and from southern Asia east to eastern India.

Importance to man. The scavenger habits of most of the members of the Hyaenidae make them of considerable benefit to the ecosystem as a whole and often directly to man. They perform the service of eliminating the remnants left by the large predators and thus fit into a niche unused by most mammals. In some parts of Africa, hyenas are allowed to roam the villages, cleaning up garbage. In some areas people place their dead out in the open to be removed by these inexpensive "natural undertakers."

Hyenas at times will attack and kill livestock and sometimes people, usually sleeping persons or young children. Being cowardly, they seldom fight if the individual puts up any type of defense. When hungry, they may raid camps and carry off various items made of leather. In attacking large livestock or antelope, they usually approach from the rear and kill their prey by tearing open the abdomen.

Form and function. Hyenas have a total of 34 teeth, fewer than the generalized number found in other carnivores. The aardwolf has fewer upper and lower premolars and molars, resulting in from 28 to 32 teeth. The incisors are unspecialized in all members of the Hyaenidae, with the third incisor being the largest. The canines are well-developed, sharp, and elongate. Premolars of hyenas have well-developed crowns, useful in bone crushing, and tend to be conical. In the aardwolf the premolars are small, widely spaced, and, along with the small molars, tend to be vestigial. The lack of development of the cheek teeth is related to its insect-food diet, which does not require a powerful dentition. The molars of the hyenas are large and powerful, aiding in bone crushing, as do the jaws, the most powerful found in any mammal of comparable size. The head tends to be massive, with a broad muzzle, not long-nosed as in the canids or short-faced as in the cats. Anal glands are present in most forms and are well developed in the aardwolf, which can eject an obnoxious fluid when threatened.

All hyenas have a characteristic stance. The front legs are longer than the hindlegs, causing the back to slope conspicuously from a high point at the shoulders to a low point at the hindquarters. Hyenas walk on their toes and run well. There are four toes on each foot, capped by blunt, nonretractile claws. *Proteles* has five toes on the front foot. Most members of the family have a well-developed mane on the neck, which may extend down the midline of the back. This mane often is more noticeable in the young hyenas, although well developed in all age classes of the aardwolf. All have large, rounded ears. These animals are grayish or brownish in colour, marked with spots (*Crocota*), stripes (*Proteles*, *Hyaena hyaena*), or uniform with bars on the legs (*Hyaena brunnea*). The tail is medium in size, bushy at the base, and pointed at the end.

Cats (family Felidae). *Natural history.* Almost all cats are strictly carnivorous and feed on small mammals and birds or on the larger herbivorous artiodactyls like deer and various types of antelope. The fishing cat (*Felis viverrina*) feeds largely on fishes and clams or snails and thus fits into a slightly different niche from that of most cats. The flatheaded cat (*Felis planiceps*) is the only felid known to feed to any extent on vegetation, with fruit and items like sweet potatoes being preferred when available. The large cats sometimes drag the kill into a tree or place it under a bush after the initial gorging. Cats live on a feast-or-famine routine, gorging themselves when a kill is made and then fasting several days between kills.

Cats have good senses of sight and hearing. Long, sensitive whiskers on the face aid the cat during the stalking of the prey by brushing against any obstacle in a path and enabling the cat to avoid making excessive noise as it stalks its prey. The easiest, most open path can be fol-

Hunting
method of
cats

lowed through vegetation, even at night. Smell apparently is not as well developed as in some other carnivores, such as the dogs. Cats are nocturnal in habits. The cheetah is an exception to this rule, hunting predominantly by day. Their large eyes are especially adapted for seeing at night, when the amount of light is low. The cat stalks its prey to a point as close as allowed by available cover, then closes the final distance by a leap or a short dash. A chase may ensue, with the predator relying on superior speed and reflexes to overtake the dodging prey, which often has greater endurance. Although the endurance of most felids is sufficient for a chase of only a few hundred metres at most, the cheetah may pursue its prey for much greater distances, up to about 5,500 metres (3.4 miles). If overtaken, the prey is thrown down and dispatched with a deep bite, usually in the neck. Almost all hunting by these animals is done alone and quietly. The major exception to the solitary habit is the lion (*Leo leo*), in which the group, or pride, may number as many as 30 individuals. Cheetahs sometimes are found travelling in small groups but stalk and chase their prey individually.

The gestation period of most smaller cats is approximately two months (50 to 68 days), and that of the large cats is three or three and a half months (88 to 113 days). Two or three kittens make up the usual litter. The jaguar (*Leo onca*) tends to have only one kitten. The domestic cat (*Felis catus*) sometimes has more than six kittens. Female cats may have from four to eight nipples. The breeding season usually is in the late winter or early spring. In smaller species those females that have early litters may produce a second litter in late summer or early fall. Some cats are capable of breeding at any time during the year (lion, tiger, and leopard). The size of the animal does not seem to determine the litter size, number of litters, or the time of the breeding season. In the larger cats, however, the age of initial breeding is greater; the females may be three or four years of age and males as much as five or six years old. Smaller cats may breed when less than a year old. Most litters are born in places seldom disturbed by man, such as in a rocky cavern, under a fallen tree, or in a dense thicket. The serval (*Felis serval*) uses an old porcupine or aardvark burrow. In most species the male does not aid in the care and raising of the young, and in fact, the female may have to guard against his attacks on the kittens. The male jaguar and ocelot (*Felis pardalis*) do help in raising the young.

Vocaliza-
tions of
cats

Cats are noted for their purring when apparently content and for their snarling, howling, or spitting when in conflict with another member of their kind. The larger forms, especially the lion, often roar, growl, or shriek. Usually, however, cats are silent. Many cats have "clawing trees," upon which they leave the marks of their claws as they stand and drag their front feet downward with the claws extended. House cats are well known for this habit, and an owner may supply an artificial tree to save damage to furniture. Whether such behaviour is for the purpose of cleaning or sharpening the claws is debatable, but the behaviour is innate; kittens raised in isolation soon begin to claw objects. Other characteristics of cats are the constant movement of the tail when stalking, the habit of washing the face, and the habit of digging a hole in which to bury the fecal material and urine.

Almost every area on earth has some type of native cat. Australia is the only continental area that has no native felids, though the family also is absent from Madagascar and a number of oceanic islands. Cats are found in most habitats, from the forest to the desert, although they are typically animals of wooded environments. Many are in danger of extinction because of their incompatibility with the activities of man or because of their combined value to zoos and the fur trade.

Importance to man. The cat, like the dog, was originally taken into the home of man as a pet. Unlike the dog, the cat has not been trained to hunt or to serve as a guardian. Cats may perform the function of "mouse catcher" in those areas where mice become a problem, but most urban house cats do not have the opportunity to perform this unless garbage and trash are allowed to accumulate. On farms, where storage of grain is common

and rodents become numerous, cats help to hold their numbers in check.

The fur of cats is sometimes in great demand, especially when high fashion calls for "fun furs," those with contrasting colours involving spots or stripes of the type found on the pelts of many cats. The demand is such that numerous kinds of cats are in danger of becoming extinct. In some regions of the world, catskins have been in great demand for many years. In parts of Africa the royal furs are usually those of a cat, as those in Eurasia are from the ermine or the sea otter. Chiefs and other high dignitaries of some native tribes use the pelt of the serval as a sleeveless cloak thrown over the shoulders.

The larger cats are strong, fierce, and extremely dangerous when hungry. Should one of these animals learn that the flesh of man is edible, it must be eliminated. In many, if not most, cases the "man-eater" is an old cat no longer able to kill native wild prey. Although tigers and leopards are most famous for man-eating activities, lions and jaguars also may become dangerous. The American lion or puma (*Felis concolor*) tends to avoid contact with man and most records of its attacking man, except when cornered, are questionable. Similarly, if a large cat learns that a ready meal is available in the form of fenced or domesticated livestock, the rancher will have problems until the predator is eliminated. As in the case of the man-eater, the problem is caused by one animal and not all the individuals should be condemned.

Cats are intelligent and can be trained. Many kinds of cats can be tamed as pets, but some, especially the larger species, really can never be trusted when they become older or when they are sexually active. Most wild animals, to be trustworthy, must be fed by bottle before their eyes open, in which case the cat apparently will regard the human as its actual mother. The domestic cat, when liberated in the wild or allowed to roam in the field, may become a more serious and damaging predator on wildlife than the native predators in the region.

Form and function. Cats have a reduced number of premolar and molar teeth; the typical dental formula includes only 30 teeth. The incisors are small and chisel-like and the canines long and pointed. The premolars are sharp, and occasionally an upper premolar may be lacking. The lower molar is elongate and sharp, the upper molar rudimentary. Because of the reduction in the number and size of the cheek teeth, a space remains between the canines and premolars in all cats except the cheetah. Felids form the most carnivorous group in the order, and the highly developed carnassial teeth reflect this specialized food habit. There is little if any specialization in the teeth for grinding or chewing. The strong masseter (chewing) muscles of the jaw restrict the amount of lateral movement and primarily permit vertical movements of the jaw, for holding the prey in a viselike grip and for slicing off pieces of meat with the carnassials. Thus meat is cut off and swallowed in relatively unchewed chunks that are broken down by strong enzymes and acids in the digestive tract. Cats use the rasplike tongue to remove the flesh from the bone.

Because of the reduction in the number of teeth, cats have a short face and a rounded, compact head. The ears tend to be short or, in some forms, tufted with hair. Most forms, except the lynxes (*Felis lynx*, *Felis rufa*), have an elongate tail that comprises about a third of the total length. The agility of a cat is evident in its anatomy. The clavicle, or collarbone, is much reduced in size. It does not connect with other bones but is buried in the muscles of the shoulder region. This type of construction allows the animal to spring on its prey without danger of a supporting or connecting structure, like the clavicle, being broken. The hindlegs are well developed, with powerful muscles that propel the animal in its spring toward or onto the prey animal. In addition to the power of the hindlegs, the animal uses strong back muscles to straighten the spinal column and provide extra force in springing and running. Most of the skin covering the body of the cat is loose, allowing for even more freedom of movement as well as providing some protection during a fight.

The legs of the jaguarundi (*Felis yagouaroundi*) are rather

Man-eaters

Anatomical
basis
of feline
agility

short, and those of the cheetah elongate. All cats, except the cheetah, have retractile claws. Such claws can be drawn back into a sheath when the animal is walking, rendering the footsteps noiseless and keeping the claws sharp; or they can be extended as an aid in pulling down and holding the prey when the animal is attacking or in protecting the cat when attacked. The sharp, strongly curved claws are also utilized when the cat climbs trees. There are five toes on the front foot and four on the rear. The first toe and its pad on the front foot are raised so that only four toes register in a track.

Cats are usually some tone of brown in colour. Some, like the American mountain lion and jaguarundi, are uniform in coloration, but most are marked with spots, stripes, or rosettes. The young of the mountain lion are spotted, which may indicate a spotted ancestor. Characteristic of all cats is a dark stripe extending laterally from the external corner of the eye. This stripe tends to exaggerate the expression on the animal's face toward one of meanness but probably serves to conceal the face by breaking the roundness of the eyes. Many of the colour patterns are repeated again and again in the different species of cats, from the largest to the smallest. Several species have black or nearly black colour phases mixed into normally coloured populations. The only cat with a well-developed mane or long hair on the back and chest, or both, is the male African lion, which is also the only terrestrial member of the order Carnivora that shows obvious sexual dimorphism. In many felids the male is larger than the female, but this sexual difference is not easy to distinguish, because the difference in size is slight.

Walrus and seals (families Otariidae, Odobenidae, and Phocidae). *Natural history.* Members of the families Otariidae (eared seals), Odobenidae (walrus), and Phocidae (earless seals), commonly called pinnipeds, are strictly carnivorous and mostly marine. A few venture up freshwater rivers, and two forms (*Pusa* species) are landlocked on inland lakes. Pinnipeds are amphibious, being aquatic as to food habits but terrestrial for mating, bearing young, and resting. The diet of pinnipeds consists mainly of fishes, cuttlefishes, octopuses, and crustaceans. Most of the fishes fed upon by seals are those which school. A few pinnipeds are specialized for feeding on krill (macroscopic pelagic plankton, primarily small shrimplike euphausiids). Seals known to feed largely on krill are the ringed seal (*Pusa hispida*), harp seal (*Pagophilus groenlandicus*), and crabeater seal (*Lobodon carcinophagus*); all take fishes as well. The crabeater seal is misnamed, as it is not known to feed on crabs. Some seals, such as the leopard seal (*Hydrurga leptonyx*), are quite predatory, feeding on penguins, other birds that land on water, and other seals. The Australian sea lion (*Neophoca cinerea*) feeds largely on penguins. Such predators take advantage of the surface light at night and probably feed by visual discrimination from below. The bearded seal (*Erignathus barbatus*) and

walrus (*Odobenus rosmarus*) are bottom feeders, consuming largely sessile organisms such as mollusks. The walrus in addition occasionally feeds on small whales, narwhals, and seals.

Generally, pinnipeds are adapted for excellent hearing, especially underwater. Some studies have even suggested the use of echolocation (reflected pulses of sound) in their movements and food-seeking activities, though conclusive evidence for this is lacking. Pinnipeds have a fair sense of smell on land, good enough to enable mothers to recognize their pups. Seals probably have an excellent sense of touch in the whiskers, which aids in locating their prey or warns the animal of a possible collision with a submerged object. Young walrus have such well-developed whiskers that these appear as a brushlike mustache. The walrus and bearded seal utilize these whiskers to filter food organisms found on the bottom and as a tactile organ. Phocid seals have patches of elongate hairs above the inside corners of the eyes.

Pinnipeds may be solitary at certain times of the year, as is the Ross seal (*Ommatophoca*), which lives alone in winter. Most, however, are usually gregarious, much more so than terrestrial carnivores. During the breeding season more than a million individuals may congregate on an island. Each male sets up a territory and gathers as many females as he can. The territory is maintained by frequent loud vocalization, a threat stance with neck outstretched, rushing at an approaching rival, and possibly through the use of olfactory odours thought to be produced in the mouth. Although the term harem is usually used to describe the large number of females in the territory of a male, this concept is misleading, since the females can move from one group to another. The groups are a result of the females staying clear of the males as the latter patrol their territories. Once the territory is established, the male remains to protect it, going as long as two months without eating (sea lions and fur seals). Eared seals (the otarids) and walrus tend to gather harems, while the earless seals (the phocids) are mostly monogamous (one male, one female) in their breeding activities. The gathering of harems forces the mating activity to be terrestrial rather than aquatic. In many, if not most, of the northern species of seals the breeding season is restricted in many ways by the environment. In most areas only four months exist with a climate favourable for the production of young.

Pinnipeds have a comparatively long gestation period, ranging from eight months in the leopard seal to 12 months in fur seals and sea lions. Breeding frequently takes place soon after the birth of the young (five days in the fur seal) or, as in the leopard seal, four months after the young are born. In most instances, copulation is on land, but some phocids mate in the water. The majority of pinnipeds exhibit delayed implantation in embryonic development. Delayed implantation allows births to occur at about the same time each year. The young are born either on land or on

Territoriality in pinnipeds

Pinniped feeding



Figure 36: Representative carnivores: pinnipeds.

ice. Usually there is but one young born, although most females have two pairs of mammary glands (some phocids have one pair). The mother goes out to sea to feed soon after the young is born. Since she must search to find her young when she returns, there has probably been a selection against the production of more than one young, at least in highly gregarious forms, for the time spent searching for the second pup would be at the expense of the first. It is evolutionarily better to produce one vigorous offspring than two weak ones. Nursing usually takes place on land. At birth the newborn pup can travel on land and swim, although young otarids are several weeks old before they develop enough blubber to keep them floating and insulated against the cold water. The young of the earless seals nurse for only three to six weeks, and the females remain near the young throughout this period. Pinniped milk is low in water content and high in fat. Such high fat content results in the rapid growth of the young, but perhaps of more importance is the low water content, of benefit to the mother in a habitat where water conservation is so important.

Female pinnipeds are sexually mature when two to eight years of age, usually at the age of three or four years. Males, at least in the polygamous species (those which gather a "harem"), may not be allowed to mate for several years after reaching sexual maturity. In these species the younger males are driven by the old males to a bachelor group located outside of the breeding group and restricted to a given area. After the breeding season, pinnipeds are largely pelagic (open-sea dwellers), travelling long distances either singly or in small groups.

A number of animals prey on seals: large sharks, killer whales, polar bears, and other seals, such as the leopard seal and walrus. Seals, especially the young, spend a good deal of time playing. They resemble the river otter in this regard.

Pinnipeds produce a number of different types of sounds. Barking by male seals may be related to social dominance and territorial defense, as the polygamous species are extremely vocal and the monogamous forms extremely quiet. The "barking" of the circus seal (the California sea lion, *Zalophus californianus*) is familiar to most people. The study of underwater communication in seals is difficult and in its infancy, but a variety of sounds have been recorded. These sounds may play a role in navigation, social behaviour, and foraging. Seals also roar or bellow, honk, chirp, bleat, grunt, or cough.

These aquatic, carnivorous animals are found throughout most of the coastal regions of the world. They are especially abundant in the cold polar waters. The walrus and phocids are well distributed throughout the northern coastal and ice-front areas. The otarids do not extend as far north, southern Alaska and Kamchatka being the most northern points, and are absent from the mid-Atlantic and North Atlantic Ocean areas. Pinnipeds are absent from the region of the Indian Ocean, southern Asia and northern Australia, the coastal regions of Central America and northern South America, and most of the coastline of central Africa. A few seals are landlocked on the Caspian Sea (*Pusa caspica*) and Lake Baikal (*Pusa sibirica*), where they subsist mainly on fish.

Importance to man. The seal is of significant importance to the survival of the Eskimos and other inhabitants of the north, who use almost every part of the animal: the meat is eaten; bones used for implements; tendons for sewing; hides for leather to make footwear, boats, shelters, bags, clothing, and ornaments; and the oil and fat for fuel, lubrication, and tanning. The ivory tusks of the walrus are carved into statues and knife handles. Seals also are taken commercially for their oil, for meat that can be used for food or fertilizer and for their hides that are used as leather.

A few seals are of importance to the fur trade. Of special interest are the true fur seals, *Callorhinus ursinus* in the northern regions of the Pacific Ocean and *Arctocephalus* species in the southern regions off the coasts of South America, Africa, and Australia. Most of the species utilized by the fur trade are members of the family Otariidae, which have dense underfur. The young of members of

the Odobenidae and Phocidae have an insulating coat of short, dense, wool-like hair at birth. Skins of these young animals are of some interest to the fur trade but are used mostly for novelties. Adult phocids have skins with stiff hair; the underfur is reduced to the point that the hide appears naked. However, some of the extremely striking phocids were in demand for furs in the late 1960s.

California sea lions are trained to entertain at circuses and zoos. Some seals become pests to commercial fisheries and may gather during salmon or herring runs, sometimes being killed in large numbers by fishermen.

Form and function. Basically, the difference in the tooth pattern of pinnipeds, in comparison to the dentition of terrestrial carnivores, reflects an adaptation toward grasping and tearing rather than chewing. Thus there is less variation in the form of the different teeth (commonly said to be homodont, or "similar teeth"), and there tends to be a reduction in the number of teeth. There are, generally, two or three pairs of upper incisors and one or two pairs of lowers. There is a pair of conical canines that form elongate tusks in the walrus. The premolars and molars are conical and so similar that they are difficult to differentiate and are usually referred to as postcanines. There are no carnassial teeth. Usually there are ten postcanine teeth each in the upper and lower jaws, but the exact number can vary so greatly, even within one species, that the total number of teeth may range from 26 to 38. In the walrus the tooth number is reduced even further, the usual number being two upper incisors, no lower incisors, four canines, and 12 postcanines, a total of only 18 teeth. This number is variable, with some individuals having as many as 24 teeth. Weddell's seals have canines and incisors adapted for cutting holes in the ice so that areas are available for breathing when the water freezes over. The upper incisors extend forward and contact the ice before the canines.

Pinnipeds have a number of aquatic adaptations that separate them from the terrestrial carnivores. The body is streamlined, allowing the animal to pass through the water with the least amount of friction. External ears are reduced or absent, also aiding in making the pattern of the body more streamlined. Adaptive features of the skull include an overlapping of the bones, a feature even more exaggerated in the whales, and a flattening of the head, an advantage during diving. The neck is thickened by increased musculature but is more flexible (at least in the eared seals) than in the terrestrial carnivores, allowing the animal to capture its prey in the water with greater ease. In most pinnipeds the face is short, with the cranial part of the skull longer than the facial region. The eyes are large, located forward, close together, and imbedded in a cushion of fat. The iris of the eye is contracted when out of water, dilated and almost circular when diving. This type of eye can accommodate rapidly when the animal moves from dim to bright light. The ears and nostrils can be closed while diving. The pelvic girdle is more nearly parallel to the vertebral column than in terrestrial forms. Such a change in the location of the pelvis enables the hindlimbs to be held in more of a trailing position and to be more efficient in propulsion through the water. The limbs are enclosed within the body skin to or beyond the elbows and knees. This enclosure also streamlines and results in a more efficient passage through the water. The extremities of the limbs are flattened to form paddles used to steer or propel the animal through its environment. Each paddle contains five toes, extending from bones in the palm and sole (metapodials), elongated to increase the surface area of the paddles. The first digit on both forefeet and hindfeet is more elongate than the other four. The tail is short and rudimentary, a condition approached by the sea otter but quite unlike the paddlelike tail of the freshwater otter. Mammary teats and external reproductive organs are so constructed that they can be withdrawn beneath the skin to aid in preserving the smooth outline of the body. The skin, adapted to a water environment, is tough and thick and contains a thick layer of subcutaneous fat or blubber. The blubber provides a source of energy that can be used during lactation or at times of fasting. Such fasting is characteristic of pinnipeds. This layer also provides the animal

Seal
barking

Adapta-
tions for
aquatic life

with more buoyancy, so that sinking in the water is less apt to occur. In addition, the blubber insulates the body so that the body temperature remains nearly constant. A pinniped living in cold water may be exposed to a surface temperature of about 0° C (32° F) and yet have a body temperature of 37° C (98.6° F), requiring an efficient system of heat production and insulation. In warm waters, and especially in warm air, the animal cannot get rid of excess body heat rapidly enough to prevent overheating. The flippers, so useful for swimming, can be used on land as radiators to get rid of excess body heat. One of the most obvious behavioural features of a colony of seals on land is the waving or extension of the naked flippers.

Physiological adaptations for diving

There are a number of physiological adaptations to life in the sea. Being air breathers, seals must be able to carry as much oxygen as possible with them in diving and to conserve what they do take. This adjustment is obtained by the increased volume of blood (very large for an animal their size) and the decreased rate in the use of oxygen by the body (decreased metabolism). The heart rate decreases in the elephant seal from about 85 beats per minute when the animal is on the surface to approximately 12 per minute during a dive. This drop in heart beat rate conserves oxygen in terms of the amount being carried in the blood. There cannot, however, be a decrease in the amount of oxygen being carried to some parts of the body. The heart and brain, for example, must have a constant and sufficient supply, or the cells in these organs will be damaged. This problem is alleviated by closing down the small arteries in the tail and flippers, while the larger arteries to the head and brain remain unchanged. The muscles function without free oxygen during submergence; the energy is obtained by partly breaking down the blood sugars and without expending the limited supply of oxygen. Most dives last less than 15 minutes, and upon return to the surface there is a greater exchange of oxygen and carbon dioxide than is found in terrestrial mammals. The amount of tidal air in one respiratory cycle may be as much as 90 percent, compared to 20 percent in man. Most dives are deep vertical dives rather than extensive horizontal dives because of the problems of orientation toward the breathing hole in the ice when horizontal directions are taken. Since there is no exchange of gases with the environment while submerged, seals have the additional problem of an accumulation of carbon dioxide. They seem to be less sensitive than terrestrial animals to increased amounts of this gas. The body temperature is lowered, a result of the decreased rate of metabolism. Seals are also adapted to withstand the great pressures during their dives. Weddell's seal (*Leptonychotes weddelli*), for example, dives as deep as 335 metres (1,100 feet), more than four times the distance that man can go in a diving suit. The longest known dive by a Weddell's seal lasted 43 minutes and 20 seconds (to 350 metres [1,150 feet]), with one elderly bull reaching a depth of 600 metres (2,000 feet) but remaining down a shorter time. Exactly how these animals can accomplish such feats is unknown. The anatomy of the seal reveals part of the answer. Many structures in the body are built to yield to such pressures, not to resist them. Some of these adaptations are smaller lungs; elimination of most air from the lungs before the dive; a thorax that permits the lungs to collapse; incomplete rings of cartilage in trachea, permitting complete collapse; driving air out of the middle ear by partly filling the ear cavity with blood; and redistribution of blood in general. The phocids are more efficient divers than the otarids, and, of the phocids, those living in the Antarctic dive the deepest.

The clavicle, or collarbone, is absent, allowing more flexibility in the shoulder. The kidneys, unlike those of most mammals, are composed of numerous lobes resembling a bunch of grapes, a condition also found in the terrestrial bear and otter.

Members of the families Otariidae and Odobenidae can walk on their limbs when on land. Phocids must move by undulations of the body as the limbs cannot be turned in a position to support the body. This does not mean that the phocids are helpless on land; the crabeater seal, a phocid, is probably the fastest moving seal on land, being difficult for even a running man to catch.

Opinion is divided as to whether the order Carnivora arose from the ancient creodonts or had a separate and independent origin from the order Insectivora. Most paleontologists now favour the latter view. The Cretaceous insectivore *Procerberus* seems to be morphologically close to the most primitive carnivores (miacids), the creodonts, and the primitive hoofed mammals (ungulates). *Procerberus* was an unspecialized rat-sized predator. The order seems to be a separate and distinct line that includes the modern Fissipedia (terrestrial forms) as well as Pinnipedia (marine and aquatic forms), while also including the more ancient and extinct family Miacidae.

Although not much is known about them, the miacids were probably arboreal forest dwellers of the tropics. Animals from this type of habitat are rarely preserved as fossils, which accounts for our incomplete knowledge of the group. These primitive carnivores may have resembled our modern weasels, with elongate bodies, short limbs, and long tails. Miacids possessed several creodont features, such as an unossified tympanic bulla (a large bubble of bone or cartilage surrounding the internal ear) and lack of fusion of the carpal (wrist) bones. Miacids differed from the creodonts in having a larger brain capacity, toe sections that were not grooved, and well-developed carnassial teeth. The carnassials, as in modern carnivores, involved the fourth upper premolar and first lower molar.

In the late Eocene and early Oligocene these primitive miacids underwent an adaptive radiation; *i.e.*, they produced diverse lines that seem to have been the beginnings of the present-day families of the Carnivora. The changes needed to convert a miacid into a modern terrestrial carnivore are minor: fusion of the separate bones of the carpus and ossification of the tympanic bulla. Carnivores that began to resemble the modern weasels, dogs, and civets were in evidence in the Oligocene, but these early carnivores were still quite similar to each other and had not differentiated to the extent of representing different families.

Two major lines within the Carnivora were distinct by the late Eocene. One, the Feloidea (or Aeluroidea), today contains cats, viverrids, and hyenas; the other, the Arctoidea (or Canioidea), contains the mustelids, dogs, bears, and raccoons. Some fossils are intermediate between the two lines, but most show definite relationships to one or the other.

Feloids had the tympanic bulla made up of an external (tympanic) and an internal (endotympanic) bone, each being a separate ossification. In many forms the claws were retractile. The transition between the miacids and viverrids (civets) was so gradual that some authorities have placed the miacids in the family Viverridae. Because of the similarities to the miacids, the viverrids are considered to form the basal stock of the feloid line. Hyenas appear to be a more recent line, having evolved from the viverrids in the early Miocene. *Ictitherium*, a mid-Miocene viverrid, had characteristics intermediate between the viverrids and hyenas, but apparently was more like the "true" hyenas. Other fossil viverrids also showed similarities to the hyenas. Although the cats seem to have evolved "suddenly," the resemblance of certain catlike viverrids (such as the Oligocene *Stenoplesictis* and the modern *Cryptoprocta* of Madagascar) to the recognized felids is regarded by some authorities as indicating the close relationship of the felids to the viverrids. In the late Eocene and early Oligocene, the felids were quite distinct, with a number of highly specialized cats already present. Of particular interest are the sabre-toothed cats such as *Hoplophonus* and the false sabre-tooth *Dimictis*, which had large upper canines fitting into a flange on the lower jaw. The Pleistocene *Smilodon*, a more recent genus, is probably the best-known sabre-toothed cat. The genus *Felis* dates back to the early Pliocene, by which time its members had features similar to those of the modern cat.

In contrast to the feloids, the arctoids almost never had retractile claws; the tympanic bulla was made up of the tympanic bone alone; and there was a canal below the bulla, through which the carotid artery passed, supplying blood to the brain. Some of the earlier miacids are difficult to separate from the mustelids, which were distinct

The extinct Miacids

Sabre-toothed cats



Figure 37: Oligocene carnivore, *Hesperocyon*, an ancestor of the dogs.

Drawing by R. Keane based on A.S. Romer, *Vertebrate Paleontology*, Copyright © 1966 by the University of Chicago Press, all rights reserved

as early as the late Eocene. They probably were one of the first arctoid lines to differentiate. By the close of the Miocene, several general types of mustelids were evident: the weasels, badgers, skunks, and otters. Another ancient arctoid line included the dogs. Oligocene canids, represented by the genera *Cynodictis* and *Hesperocyon*, resembled elongate, short-legged weasels and civets. Both genera, although considered to be primitive dogs, were generalized enough to be ancestral to many carnivores. Genera that were more typically canid appeared in the early Oligocene. These were running animals, with four well-developed toes and reduced hindtoes (hallux and pollex), each ending in a blunt claw. More than 50 genera of doglike forms are known, but the main line of evolution seems to have passed through forms like *Temnocyon* (upper Oligocene) and *Cynodesmus* (lower Miocene) to the modern genera *Canis* (first seen in the lower Pliocene deposits), *Vulpes* (perhaps in the upper Miocene), and *Urocyon* (Pleistocene). Two lines of extinct canids are the hyena-like dogs and the bear dogs, neither of which has any relationship to the living hyenas or bears. Because of the close relationships of the bears, dogs, and raccoons, various taxonomists have disagreed on the placement of intermediate genera. For example, *Dinocyon*, *Hemicyon*, and *Cephalogale* were once considered to have been canids but are now placed with the bears and regarded as perhaps part of the ancestral lineage that diverged from the canid line in the late Oligocene or early Miocene. *Phlaocyon* and *Aleocon* were once considered to be primitive procyonids because of their teeth, but the skull and ear structures are more reminiscent of the dogs. A recent

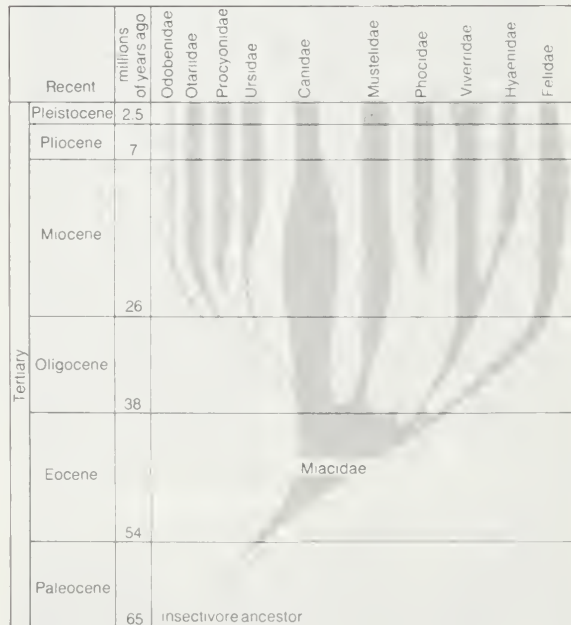


Figure 38: Phyletic dendrogram of Carnivora.

form, the giant panda (*Ailuropoda*), considered by some to be a bear and by others to be a procyonid, probably diverged from primitive ursid stock when the ancestor of the bear and raccoon lines still had features common to both. Procyonids (raccoons) are poorly represented in the fossil record, perhaps because of their arboreal habits. They appear to have diverged from the canid line in the late Oligocene.

The preponderance of carnivorous and carrion-eating fishes and mammals in the open sea, along with the buoyant characteristics of the bodies of the pinnipeds, results in the almost immediate destruction of any seal that dies. Thus, dead pinnipeds seldom settle intact and become buried (the first steps toward becoming a fossil) and are poorly represented in the fossil record. Relationships must instead be inferred through studies of recent forms. Although there is no fossil record of pinnipeds prior to the Miocene, it is probable that the ancestors of the pinnipeds evolved during the adaptive radiation of the miacids in the Eocene.

There are two rather distinct evolutionary lines in the animals called pinnipeds. One includes the two rather closely related families Otariidae and Odobenidae. This line may have become distinct in the early Miocene or late Oligocene, evolving from a form closely related to the canid ancestor of the bears. The other line led to the modern phocids, which are quite different from the otarid-odobenid group. Their affinity to any other group of the Carnivora has not been shown, although the mustelids have more frequently been suggested as related to the phocids than any other group. *Semantor*, an aquatic carnivore from the lower Pliocene, was once considered to belong to a separate pinniped family but is now considered a mustelid. The separation of the phocid and otarid-odobenid lines must have occurred at an extremely early date, if the two lines ever had a common pinniped ancestor. Some paleontologists postulate a separation in the early Miocene or late Oligocene. Possibly the two lines never evolved from a single aquatic group but had separate origins within the miacids, sometime in the late or even mid-Eocene. Or perhaps the otarids stemmed from an ancestral ursid type in the early Oligocene. Such an early separation, or lack of any real ancestral association, could have produced the dissimilarity of the two types of pinnipeds, as well as the difficulty in associating the phocids with any other family of the Carnivora.

Evolution of pinnipeds

CLASSIFICATION

Distinguishing taxonomic features. If one groups the seals with the terrestrial carnivores, the number of characteristics common to all members of the order Carnivora is small.

The characteristics used to separate the Carnivora from other mammalian orders and to define the subdivisions of the Carnivora are primarily structural. Of great importance are certain features of the skull (the type of jaw articulation and the shape of the paroccipital process), of the feet (the number of toes, lack of opposability of the hindtoe, type of claws, and fusion of the scaphoid and lunar bones), and of the teeth (both the overall tooth pattern and the shape of the individual tooth). The dentition is especially important in determining the relationships of fossil forms. Also useful in the taxonomy of modern carnivores are the convolutions (irregular ridges) around the lateral or sylvian fissure of the brain, the relative weights of the adrenal and thyroid glands, the type of uterus and placenta, and the position of the nipples.

Annotated classification. Although the terrestrial and marine carnivores are grouped together in the following classification, many authors separate them into two orders. In this article the split is indicated by erecting two suborders, the Fissipedia and Pinnipedia, based largely on the degree of specialization undergone by the members as they evolved and adapted to radically different environments. Basically the taxonomy presented combines the classifications by G.G. Simpson and A.S. Romer, North American paleontologists.

Numerical dental formulas are used below to indicate the number of each type of tooth in the right or left half

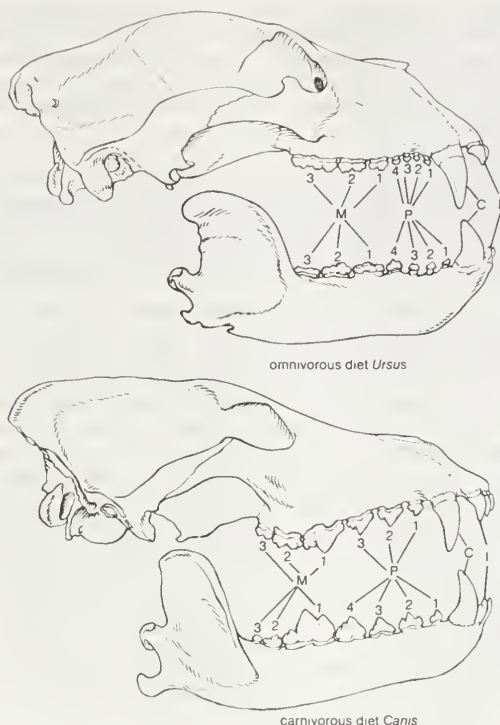


Figure 39: Tooth patterns in omnivorous and carnivorous members of Carnivora. Side view of lower mandible and cranium.

From Hall and Kelson, *The Mammals of North America*, vol. II, copyright © 1959 The Ronald Press Company, New York.

of the upper and lower jaws and the total number of teeth in the mouth. Four tooth types are indicated in the following manner:

Upper jaw: incisors . canines . premolars . molars = total number of teeth
Lower jaw: incisors . canines . premolars . molars

ORDER CARNIVORA

Medium- to large-sized terrestrial or aquatic meat eaters distributed worldwide; composed of 10 extant families, about 115 living genera, and 274 living species.

Suborder Fissipedia

Terrestrial (except otters and polar bear); usually 3 pairs of incisors above and below, adapted for biting off flesh; premolars vary in number from the usual 4 on each half of each jaw (4/4) to 4 above and 2 below (4/2); molars vary from 2/3 to 1/2; carnassial teeth well developed and formed by the 4th upper premolar and 1st lower molar; usual dental formula: $\frac{3 \cdot 1 \cdot 4 \cdot 2}{3 \cdot 1 \cdot 4 \cdot 3} = 42$ teeth; all-inclusive dental formula: $\frac{3 \cdot 1 \cdot (2-4) \cdot (1-4)}{(2-3) \cdot 1 \cdot (1-4) \cdot (1-5)} = 28-50$ teeth; facial part of skull usually elongate; tympanic bullae large; alisphenoid canal (a perforation in the skull for a branch of the carotid artery) usually present; lacrimal duct present; pelvis at distinct angle from vertebral column; femur (thigh bone) straight, moderately slender, with small head; fibula (minor lower bone of hindlimb) much smaller than tibia (shinbone); 1st metatarsal and bones of phalanges with epiphyses only at proximal ends; metapodials never elongate; vertebrae with interlocking processes; external ears (pinnae) well developed; intestine short (6 to 8 times length of body in dog); kidneys simple, not lobulate except in otters and young bears.

Superfamily Miacoidea

All extinct generalized carnivores with carnassial teeth. Middle Paleocene to upper Eocene. There is one family, Miacidae, about 13 genera.

Superfamily Feloidae (cats, civets, and hyenas)

Tympanic bullae divided, composed of external tympanic and internal endotympanic bones; paroccipital process in close contact with bullae; most have retractile or semiretractile claws; cecum small; baculum (bone of penis) present or absent; usually polyestrous; upper Eocene to Recent; 42 living genera and 115 species.

Family Viverridae (civets and mongooses). Upper Eocene to Recent. Distributed throughout most of the southern regions of the Old World (Asia, Europe, Malay Peninsula, most of Africa, Madagascar, and associated southern oceanic islands) but most abundant in the India-Malay region; introduced

on New Zealand, Hawaii, and the West Indies. Small to medium sized (0.7 to 14 kg; 1.5 to 31 lb). Well-developed carnassials; upper carnassials usually lack an anterior lobe and the lower carnassials have a well-developed talonid process; middle incisor on each side located farther inward than the outer 2; incisors broad; 4 small premolars, with 1st minute or absent; molars 2/2 or in some cases 1/2 (*Prionodon*), 2/1 (*Salanoia*), or 1/1 (*Cryptoprocta*); molars large, especially the 1st; dental formula usually $\frac{3 \cdot 1 \cdot (3-4) \cdot 2}{3 \cdot 1 \cdot (3-4) \cdot 2} = 36-40$ teeth.

Alisphenoid canal usually present (absent in *Eupleres* and most members of subfamily Galidiinae); entepicondylar foramen of humerus usually present; auditory bulla externally constricted and divided by septum. Baculum well developed. Intestinal cecum small, sometimes absent. Elongate body, short legs; digitigrade or semiplantigrade gait; 5 toes usually present, with semiretractile claws (retractile in *Cryptoprocta*); pollex and hallux located on side of foot above other toes. Ears small and rounded; tail usually long (12-90 cm; 4.7-35 in.) and bushy or tapering to a point from a bushy base. Scent glands often well developed, producing secretion called "civet." Usually uniform in coloration but some spotted or striped.

Subfamily Viverrinae (true civets, genets, and linsangs). Mainly African; approximately 7 living genera and 15 species.

Subfamily Paradoxurinae (palm civets). Mainly Asian, 6 living genera and 8 species.

Subfamily Hemigalinae (banded palm civets, web-footed civets). Asia or Madagascar, 5 living genera and 7 species.

Subfamily Galidiinae (striped and ring-tailed mongooses). Restricted to Madagascar, 4 living genera and 7 species.

Subfamily Herpestinae (true mongooses and meerkats). All African, 13 living genera and approximately 37 species.

Subfamily Cryptoproctinae (*fossa*). Monotypic, *Cryptoprocta ferox*, found only in Madagascar.

Subfamily Stenoplesictinae. Fossil. Upper Oligocene to middle Miocene of Europe.

Family Hyaenidae (hyenas and aardwolves). Middle Miocene to Recent. Found today throughout most of Africa and southern Asia (Turkey to eastern India). Medium-sized (10 to 80 kg; 22 to 175 lb) carrion feeders; teeth specialized for crushing bones; incisors unspecialized, outer (3rd) larger than other 2; canines powerful; premolars strong and conical; molars strong and large; carnassials well developed (hyenas); long, slender canines in aardwolf (*Proteles*) with small and widely spaced premolars, small molars, and no development of carnassials.

Dental formula $\frac{3 \cdot 1 \cdot 4 \cdot 1}{3 \cdot 1 \cdot 3 \cdot 1} = 34$ teeth for hyenas; premolars 3/(2-1), molars 1/(1-2) = 28-32 for aardwolf. Alisphenoid canal absent; auditory bullae partially divided by a rudimentary septum in hyenas, divided in aardwolf; front legs longer than rear legs; digitigrade; 4 toes (5 on forefoot in aardwolf); claws blunt and nonretractile; entepicondylar foramen absent; baculum absent; more thoracic vertebrae (15) than any other feloid; tail of medium length (20 to 33 cm) and bushy. Preanal scent glands absent. Cecum 15 to 23 cm (6 to 9 in.) long. Short mane usually present. Ears large and rounded.

Subfamily Proteolinae (aardwolf). African, 1 species, *Proteles cristatus*.

Subfamily Hyaeninae (hyenas). Africa and southwest Asia, 2 living genera, 3 species.

Subfamily Ictitheriinae. Fossil. Miocene-Pliocene, 1 Old World genus.

Family Felidae (cats). Upper Eocene to Recent. Now worldwide except for Antarctica, Australia, Madagascar, and some oceanic islands. Size medium to large (2.5 to 275 kg; 5.5 to 600 lb). Highly carnivorous; well-developed carnassials. Total number of teeth reduced from primitive carnivore condition; dental formula typically $\frac{3 \cdot 1 \cdot 3 \cdot 1}{3 \cdot 1 \cdot 2 \cdot 1} = 30$ teeth; alisphenoid canal absent; auditory bullae not constricted externally; posterior palatine foramina on maxillopalatine suture instead of on maxilla. Elongate body with relatively short legs; digitigrade; 5 toes on front, 4 on rear feet, capped with sharp, curved, retractile claws. Entepicondylar foramen present; baculum absent or vestigial; cecum small; tail short (lynx) to long (up to a third of total length). Preanal scent glands absent; tongue covered with horny papillae.

Subfamily Felinae (lynxes, typical cats, lions, tigers, leopards). Worldwide except Antarctica, Australia, Madagascar, and most oceanic islands, 2 living genera, 35 species.

Subfamily Acinonychiinae (cheetah). Africa and southwest Asia, 1 species, *Acinonyx jubatus*.

Subfamilies Proailurinae, Nimravinae, Machairondontinae, Hyainailourinae. All fossil. About 39 genera.

Superfamily Canioidea (dogs, bears, raccoons, and weasels)

Tympanic bulla formed from tympanic bone only; not divided into 2 chambers; long canal beneath bulla; paroccipital process prominent and independent of bullae. Claws nonretractile.

tile. Cecum present or absent; baculum well developed; usually monestrous (1 litter of young per year). Eocene to Recent; 53 living genera and 133 species.

Family Canidae (dogs, foxes, jackals). Eocene to Recent. Now worldwide except Antarctica and most oceanic islands. Size medium (1.5 to 80 kg; 3.3 to 175 lb). Carnivorous or omnivorous; well-developed carnassials; dental formula typically $\frac{3 \cdot 1 \cdot 4 \cdot 2}{3 \cdot 1 \cdot 4 \cdot 3} = 42$; bush dog (*Speothos*) with (1-2)/2 molars (1 upper molar is sometimes minute or absent) = 38 to 40 teeth; Cape fox (*Otocyon*) with (3-4)/(4-5) molars = 46 to 50 teeth. Facial part of skull elongate; alisphenoid canal present; legs long and semirigid; digitigrade with 5 toes on front (1 dew claw) and 4 on rear foot; claws blunt and nearly straight. Entepicondylar foramen absent; baculum well developed and grooved. Cecum always present, short and simple or long and characteristically folded. Tail relatively long (up to about a third of total length) and bushy; scent glands often on dorsal part of base of tail. Ears pointed, erect, and large.

Subfamily Caninae (dogs, jackals, and foxes). Worldwide except Antarctica and most oceanic islands; dingo probably introduced on Australia in prehistoric times; wild dogs of New Guinea probably derived from domestic dogs. 11 living genera and approximately 37 species.

Subfamily Simocyoninae (dhole, hunting dog, bush dog). Asia, Africa south of the Sahara, and northern South America. 3 living monotypic genera: *Cuon alpinus*, *Lycaon pictus*, *Speothos venaticus*.

Subfamily Otocyoninae (big-eared fox). Eastern and southern Africa; monotypic, *Otocyon megalotis*.

Subfamilies Amphicyonodontinae, Amphicyoninae, Borophaginae. Fossil only. Total genera about 34.

Family Ursidae (bear and giant panda). Middle Oligocene to Recent. Restricted to larger land masses of world, ice floes of Arctic, and in South America to Andes; bear of Atlas Mountains of northwest Africa now extinct. Size medium to large (27 to 780 kg; 60 to 1700 lb); usually omnivorous; carnassials poorly developed; dental formula variable, but typically $\frac{3 \cdot 1 \cdot 4 \cdot 2}{3 \cdot 1 \cdot 4 \cdot 3} = 42$ teeth. Skull elongate; alisphenoid canal present. Legs short, large, and powerful; plantigrade; 5 toes on each foot, capped by elongate, powerful, nonretractile claws. Entepicondylar foramen of humerus usually absent; lacrimal bone of skull, clavicle, and cecum absent; tail short, practically absent; ears small and rounded; lips mobile and free from gum. 6 genera and 8 living species.

Family Procyonidae (raccoons, coatis, lesser panda). Upper Eocene (probably) to Recent. Fossils through Pleistocene in Europe. Except for the lesser panda (*Ailurus*) of the southeastern Himalaya region of Asia, family is restricted to northern South America, Central America, and southern North America. Size medium (0.8 to 22 kg; 1.75 to 48.4 lb); omnivorous; carnassials poorly developed; dental formula usually $\frac{3 \cdot 1 \cdot 4 \cdot 2}{3 \cdot 1 \cdot 4 \cdot 2} = 40$, with some variation in number of premolars (kinkajou 3/3, lesser panda 3/4); alisphenoid canal absent except in *Ailurus*; legs of medium length; semiplantigrade to plantigrade; 5 flexible toes on each foot capped by nonretractile or semiretractile claws; entepicondylar foramen sometimes present; baculum well developed and, except for *Ailurus*, bilobed and hooked at distal end; cecum absent; clavicle vestigial; tail long (20 to 70 cm; 7.9 to 27.6 in.), often ringed, prehensile in *Potos*; ears small to medium in size and rounded. 7 living genera and approximately 14 species.

Family Mustelidae (weasels, badgers, skunks, otters). Eocene to Recent. Now worldwide except Madagascar, Australia, and most oceanic islands; introduced into New Zealand; more common in Northern Hemisphere. Size from smallest of carnivores (35 gm; 1.2 oz) to medium; mainly carnivorous; carnassials usually well developed; dental formula usually $\frac{3 \cdot 1 \cdot 3 \cdot 1}{3 \cdot 1 \cdot 3 \cdot 2} = 34$ teeth but may range from 30 to 38 teeth, with variation in incisors (3/2 in sea otter), molars (1/1 in ratel and African striped weasel), but especially premolars (2/3 in the African striped weasel, 4/2 in the river otter, 4/4 in American badger); alisphenoid canal absent; legs short in relation to elongate body; plantigrade to digitigrade; 5 toes on each foot, capped with nonretractile claws; entepicondylar foramen present or absent; tail usually long (up to 1/3 of total length); ears small and rounded; anal scent glands well developed.

Subfamily Mustelinae (weasels and Old World polecats). Subfamily with widest distribution and most kinds. There are 11 genera and 33 recent species.

Subfamily Mellivorinae (ratel or honey badger). African and southern Asia. There is 1 species, *Mellivora capensis*.

Subfamily Melinae (badgers). North America, Europe, Asia, with most forms found in Southeast Asia. There are 6 genera and 8 recent species.

Subfamily Mephitinae (skunks). New World only. There are 3 living genera and 11 species.

Subfamily Lutrinae (otters). Worldwide. There are 5 living genera and 19 species.

Subfamily Leptarcinae. Fossil. Miocene and Pliocene. New World. 3 genera.

Suborder Pinnipedia

Aquatic; teeth nearly homodont (uniform) and adapted for grasping and tearing, not chewing; 2 pairs of incisors below; postcanines (premolars and molars) similar, never more than 2-rooted, and usually 5 on each jaw (varies from 3 to 7); carnassials absent; usual dental formula $\frac{3 \cdot 1 \cdot 5 \text{ (postcanines)}}{2 \cdot 1 \cdot 5 \text{ (postcanines)}} = 34$ teeth; all inclusive dental formula $\frac{1-3 \cdot 1 \cdot (3-7)}{0-2 \cdot 1 \cdot (3-6)} = 18-38$ teeth; cranial part of skull longer in proportion to facial part than in most fissipeds; face profusely innervated by extremely large trigeminal nerve; brain more spherical and with greater development of convolutions than fissipeds. Tympanic bullae small in some; alisphenoid canal present or absent; lacrimal duct absent; basioccipital and sphenoid relatively large as compared with fissipeds; pelvis small and nearly parallel to vertebral column; femur broad, flattened, with globular head on short neck; fibula almost as large as tibia; metatarsals and all bones of phalanges with epiphyses at both ends of shaft; metapodials elongate; 5 well-developed digits on each limb; 1 digit on manus and 1st and 5th digits of pes usually longer than other digits; feet fully webbed; nails small to absent; baculum massive. Fusiform-shaped body; external ears (pinnae) reduced or absent; tail short or vestigial (5 to 20 cm [2 to 8 in.] long) with 8 to 15 caudal vertebrae. Intestine long; cecum short; kidneys lobulate. Thick layer of subcutaneous fat (blubber). One large, precocial young.

Family Otariidae (eared seals). Lower Miocene to Recent. Coastlines of Pacific, South Atlantic, and Indian oceans. Hindlimbs useful for support on land; nails well developed on middle 3 digits, rudimentary on outer ones. Small external pinnae of ears. 2 pairs of lower incisors; upper incisors of 1st pair notched transversely; usually 20 to 22 postcanine teeth; dental formula $\frac{3 \cdot 1 \cdot 4 \cdot (1-3)}{2 \cdot 1 \cdot 4 \cdot 1} = 34$ to 38 teeth. Males are 2 to 4 times as large as females (males 120 to 1,000 kg [265 to 2,200 lb]; females 60 to 270 kg [130 to 600 lb]); testes suspended in distinct external scrotum; small but distinct tail; obvious neck; skin whitish or light gray.

Subfamily Otariinae (sea lions). Coastal regions of northern and eastern Pacific Ocean, southern Australia and southwestern New Zealand; coastal regions off eastern coast of South America. There are 4 living genera and 5 species.

Subfamily Arctocephalinae (Pacific and southern fur seals). Southern coastal areas of all southern continents; north Pacific region. There are 2 living genera, 7 species.

Family Odobenidae (walrus). Upper Miocene to Recent. Circumpolar holarctic distribution, seldom south of about 58° north latitude. Hindlimbs useful for support on land. External pinnae absent; postorbital process absent; alisphenoid canal present. Lower incisors absent in adult; upper canines of both sexes form tusks; usually 12 postcanine teeth; dental formula $\frac{(1-2) \cdot 1 \cdot (3-4) \cdot 0}{0 \cdot 1 \cdot (3-4) \cdot 0} = 18-24$ teeth; males larger (up to 1,270 kg; 2,800 lb) than females (up to 860 kg; 1,900 lb). Testes abdominal (internal). No free tail. Skin whitish or light gray. Monotypic, *Odobenus rosmarus*.

Family Phocidae (earless seals). Middle Miocene of Europe and North America; Pliocene of Asia; lower Pleistocene of Africa; Recent along northern and southern coastal areas, sporadic in tropical and southern regions. Hindlimbs useless on land (cannot be placed forward under body); nails equally well developed on all digits; external pinnae absent; postorbital process rudimentary or absent; alisphenoid canal absent; 2 to 4 lower incisors; usually 16 to 20 postcanine teeth; dental formula $\frac{(3-2) \cdot 1 \cdot 4 \cdot (0-2)}{(1-2) \cdot 1 \cdot 4 \cdot (0-2)} = 26$ to 36 teeth; males may be much larger than females (elephant seal, hooded seal), slightly larger (some members of subfamily Phocinae), equal in size (some Phocinae and some Monachinae), or females may be larger than males (some Monachinae); olfactory apparatus reduced; testes abdominal; stubby tail; no visible neck; skin brown or black.

Subfamily Phocinae (northern seals). Northern temperate to Arctic waters; 7 living genera and 9 species.

Subfamily Monachinae (southern seals). Tropical and Southern Hemisphere waters; 6 living genera and 9 species.

Critical appraisal. The taxonomy of the major categories of major groups placed in the Carnivora has been in a state of flux for at least 100 years, and these categories do not seem to be stabilizing, even today. Most

Divergence
of seals

mammalogists at present regard the seals and terrestrial carnivores as belonging to different orders, the Pinnipedia and Carnivora. There are, in reality, only a few features common to the seals and their terrestrial relatives because of the extensive and numerous adaptations the aquatic forms have undergone to make them efficient carnivores of the sea. Mammalogists who have studied seals intensively now realize that there is no anatomical structure unmodified by the extensive aquatic adaptations; every organ and tissue examined has been found to be different in some way from its counterpart in terrestrial forms. In many, if not most, ways there is more difference between seals and terrestrial carnivores than between the bats and the insectivores, and few mammalogists would place the bats and insectivores in the same order. Other mammalogists, tending toward conservative taxonomy, feel the relationship of the terrestrial and aquatic carnivores can be best expressed by retaining them in two suborders, the Fissipedia and Pinnipedia, of the single order Carnivora. The conservative attitude is retained in this article.

Of the 10 living families recognized in the Carnivora, two have separated from their lines most recently and are most easily associated with other existing families: the Odobenidae with the Otariidae and the Hyainidae with the Viverridae. The family Procyonidae has perhaps the greatest number of taxonomic problems. It appears to be a collection of primitive forms whose relationships are not yet clear. Most of the living genera have at one time been placed in distinct families or subfamilies (Nasuidae, Bassariscidae, Ailuridae, Bassaricyonidae, Procyoninae, Potosinae). In most other families there are also genera sufficiently different from other members of their family that they have been given separate families (the canids *Lycan* and *Otocyon*; the giant panda *Ailuropoda*; the otters; the viverrid fossa *Cryptoprocta*; and the aardwolf *Proteles*).

The arrangement of the 10 families into two distinct superfamilies, Canoidea and Feloidea (or Aeluroidea), appears to be a natural arrangement dating back to the works of W.H. Flower and H. Winge in the late 1800s. In the Canoidea, as revealed by studies in comparative anatomy and the fossil record, the families Canidae, Ursidae, and Procyonidae seem to be most closely related. Also placed in the Canoidea is the family Mustelidae, although some of the more primitive members show resemblances to the primitive viverrids as well as the canids. In the Feloidea, the families Viverridae and Hyainidae seem most closely related, the Felidae being the most aberrant.

Those families that contain rather diverse lines have been divided into subfamilies; the number of subfamilies in each family indicating the amount of evolutionary divergence that has occurred. The groups that have probably been distinct the greatest length of time have the most subfamilies; the Viverridae with six and the Mustelidae with five. Within the viverrids some of the subfamilies are distinct, a few less distinct. Most distinct are the subfamilies Herpestinae and Cryptoproctinae. The Cryptoproctinae seems to be an ancient group with the genus *Cryptoprocta* retaining characteristics of the felids as well as the viverrids, so much so that some taxonomists have placed *Cryptoprocta* in the felids, although most mammalogists regard them to be more similar to the viverrids. The remaining four subfamilies have been combined and split in a number of ways. Usually retained, along with the two above subfamilies, are the subfamilies Viverrinae and Galidiinae. Each of the five mustelid subfamilies seems to be distinct from the others. Some authorities feel that the skunks and badgers may be close enough to place them in the same subfamily. Of the five subfamilies, the most diverse is the Mustelinae, with 11 distinct genera.

What much of the discussion above means, in terms of critical taxonomic appraisal, is that a system of classification is in some ways an artificial system set up for the convenience of man. Ideally, the system reflects real evolutionary relationships, but these must be inferred from a scanty fossil record and from comparisons of modern species. Since there are differences of opinion among specialists as to which taxonomic characters should be given the most weight, there are certain to be alternate classifi-

cations, the acceptability of which depends on new information continually being discovered. Just as the animals have evolved, so does the taxonomic system. (H.J.S.)

Cetacea (whales, dolphins, porpoises)

The order Cetacea consists of a group of primarily marine mammals occurring throughout the seas of the world and in certain tropical rivers and lakes. The name whale is often applied by scientists as a general name for the larger species; most of the smaller members of the order are called dolphins or porpoises. (The traditional grouping of these mammals in the single order Cetacea is not universally accepted. This disagreement is reflected in varying systems of classification [see below, *Annotated classification and Critical appraisal*]).

General features. The order Cetacea includes three distinct suborders: Archaeoceti, an extinct group of toothed whales; Odontoceti, the modern toothed whales; and Mysticeti, the baleen, or whalebone, whales. In addition to skull characteristics, the three groups are distinguishable by dentition: the archaeocetes were heterodont (*i.e.*, the anterior teeth were quite different in shape from the posterior ones), and milk teeth in the young were replaced by permanent teeth in the adult, as in most mammals; the odontocetes are mostly homodont (teeth of uniform shape), and most lack milk teeth; adult baleen whales lack teeth, and the animal feeds by straining minute organisms from the water by the baleen plates (see below). The odontocetes include about 70 species—the porpoises, dolphins, pilot whales, or “potheads,” false killer whales, beaked whales, the sperm whale, and a few lesser known forms. They range in length from about 1.3 metres (4.3 feet) in the smaller porpoises to about 19 metres (62 feet) in the male sperm whale. The maximum length of the 10 species of baleen whales ranges from six metres (about 20 feet), in the case of the pygmy right whale (*Caperea marginata*), to about 33.6 metres (110 feet) in the largest recorded specimen of the blue whale (*Balaenoptera musculus*). A great range of adult weight is encompassed by the Cetacea, from about 45 kilograms (100 pounds) in some small porpoises to about 136,000 kilograms (150 tons) in the blue

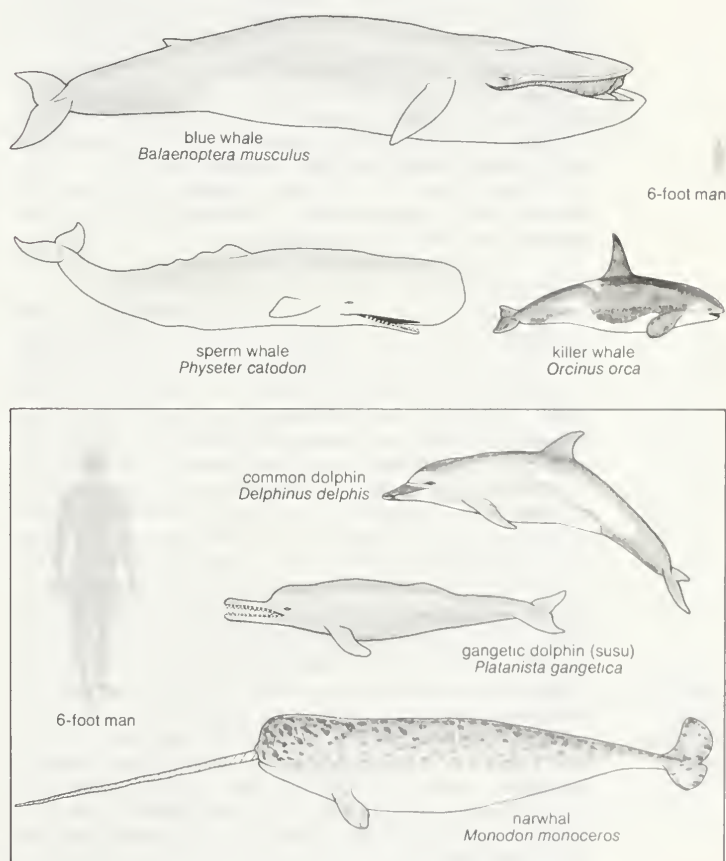


Figure 40: Body plans of representative cetaceans.

Breathing
through
blowholes

whale. The body of modern whales is fusiform (tapered); the tail, which ends in a horizontal blade (the flukes), is used for propulsion. The flukes are unsupported by bone. There is usually a dorsal fin on the back, also unsupported by bone. Breathing is accomplished rapidly at the water surface through paired (in mysticetes) or single (in odontocetes) blowholes, which are generally located on top of the head, some distance behind the tip of the snout. The exhalation produces the familiar spout, through the sudden condensation of water contained in the breath. Following inhalation, the animal holds its breath for variable periods, as it swims beneath the surface. Dives of smaller cetaceans last a few minutes, but those of some larger species may last an hour and perhaps longer.

The cetaceans are an ancient group, having split off from other mammals early in the history of the Mammalia. In spite of their many and profound adaptations for aquatic life, cetaceans are typical mammals, developing their young internally and giving birth at sea, nursing the young with milk, and possessing internal heat control mechanisms, or warm-bloodedness. Although hair is largely absent in most adult whales, the evidence from fetal or newborn animals indicates that the lack of hair is a secondary loss (*i.e.*, that hair was present in the ancestors of whales).

Importance to man. Many of the more abundant whales and porpoises are commercially important, their meat being used as food for animals and (especially in Japan) for humans and their oil for industrial lubrication and for conversion into soaps and fatty acids, which are used in cosmetics and detergents.

About two dozen countries have engaged in commercial whaling at one time or other, but by the early 1970s only Japan and the Soviet Union maintained whaling fleets in open ocean, the fleets of other countries having ceased operations as a result of declining markets for whale products and drastically reduced populations of the commercially valuable whales. Several other countries have continued to take whales from shore stations. Limitations on catches of the larger cetaceans, imposed by the International Whaling Commission (IWC), were too late to prevent the near extinction through overexploitation of the blue whale by the late 1950s. Populations of the fin whale (*Balaenoptera physalus*) have declined sharply under heavy fishing pressure. The decline in these species has led to increased pressure on the sei whale (*B. borealis*) and the sperm whale (*Physeter catodon*), by 1970. Effective regulation of commercial whaling in international waters to a level that would allow a sustained yield without depleting whale populations has been handicapped by the inability of the IWC to impose obligatory controls on its members. In 1971 the United States declared the commercially exploited whales to be endangered species and prohibited by law the importation of all whale products.

NATURAL HISTORY

Behaviour. *Social behaviour.* Most toothed whales spend their entire lives in tightly organized schools, which may range in numbers from a few animals to 1,000 or more individuals. Baleen whales are more often found singly, although small schools, or pods, occur on the breeding grounds, and whales may congregate in feeding areas in considerable numbers. Odontocete schools, at least, may be quite complex in structure, involving family groups, groups segregated by age and sex, and even schools composed of two or more species. Bottle-nosed dolphins (*Tursiops truncatus*), for example, are often found associated with schools of pilot whales (*Globicephala* species) in some areas and move with them, apparently because the food-finding capability of the large whale school is greater than that of the smaller dolphin group. Schools of toothed whales show a variety of swimming patterns, including movements as a rank or front, single file, or a more formless pattern. Within schools, subgroups usually swim, feed, and dive independently, except when the school is frightened or travelling rapidly; then members may pack together and leap and dive more uniformly. Groups of mothers with young are usually found together near the centre of a school.

Studies of captive odontocete schools, especially of bottle-

nosed dolphins, have shown that mother-young relations may persist for several years, with the young returning to its mother in times of stress even when full grown. Other nonrelated females may also be part of such groups and may assist during birth and in defense against males. Adult males tend to segregate within the school, except during the breeding season; juvenile males may gather in small groups. Dominance hierarchies exist in which both males and females exert leadership, though dominance is strongest in large adult males. Aggression by the leader is used to establish such relationships and may take the form of biting, hitting the offender with the tail, clapping the jaws intimidatingly, or ramming with the snout. Juvenile males of many species are especially apt to be heavily scarred with tooth marks. Biting or raking with the teeth is also a marked part of sexual activity, and females are therefore apt to be scarred, especially in the genital area.

Play. Play, although common in odontocetes, is almost unknown in mysticetes. In porpoises, much play is sexual in nature, consisting of rubbing or prodding the genitals with the snout or fins and flippers, usually accompanied by sound production of a complicated character. Other forms of play consist of balancing floating objects, such as kelp, sticks, or feathers, on the fins and flukes, or holding objects in the mouth. Younger killer whales and sperm whales have been seen pushing small logs or planks; food items are also used. Fish may be thrown into the air or proffered to other fish lurking in crevices, apparently to lure them into the open. Underwater observations at sea have revealed porpoise groups swimming and diving together in complicated ballet-like movements, riding ocean swells or breaking surf, in precise formation swimming. In pilot whales, playlike swimming patterns have been observed, the animals resting vertically, head down, with a metre or so of their tails in the air, or swimming on their backs.

Epimeletic behaviour. Behaviour in which one animal assists another in trouble (epimeletic behaviour) is marked in whales, having been noted in both mysticetes and odontocetes. It consists of "standing by," in which school members refuse to leave a wounded animal or actually support a sick animal, sometimes one of a different species. An extreme form sometimes occurs in porpoises: a mother may support the body of her stillborn calf until it literally rots away. Most whales and porpoises appear to be solicitous of their young, and there are several reports of babies being shielded by the mother or taken on a pectoral flipper and rolled part way out of water.

Fright. Fright is typically shown by a bunching together of school members and rapid swimming, a reaction that has been used by fishermen to capture certain species. The school is frightened by noise production underwater—banging on pipes or metal bars or clacking two rocks together. The animals are run until their shortage of breath prevents deep dives and underwater evasion, and eventually they are herded ashore or into nets. Porpoises sometimes signal moderate levels of fright by slapping the water with the tail flukes.

Feeding. Two general modes of feeding exist in cetaceans. Mysticetes are strainers who envelop large amounts of water, forcing it through horny fibrous mats of baleen (slatlike horny plates) that hang from the roof of the mouth and swallowing the food that remains on the inner surface of the baleen. The odontocetes actively pursue and capture swimming prey, especially fish and squid. Two types of feeding have been described in baleen whales: the skimming mode, in which the whale swims forward with the mouth slightly open, filtering a continuous column of water as it moves; and the gulping mode, in which the animal engulfs a mouthful of water and forces it out past the baleen before taking another mouthful.

Life history. *Reproduction.* Although all cetaceans are seasonal breeders, this period may be quite long. Adult male porpoises tend to accompany the females who are receptive. Impregnation occurs during the long period from spring to fall and, in both the large baleen whales and the porpoises, birth takes place about 11 to 12 months afterward. The sperm whale is an exception, gestation taking about 16 months. Growth is fairly rapid in all cetaceans,

Causes of
frightSwimming
patterns
of toothed
whales

and the birth size is large; a 1.6-metre (5.2-foot) harbour, or common, porpoise (*Phocoena phocoena*) is known to have given birth to a calf 0.6 metre (two feet) long, and a 30-metre (100-foot) blue whale was found to contain a full-term fetus 7.54 metres (24.7 feet) long.

Birth of whale calves takes place underwater; usually, the tail emerges first. The newborn porpoise is often pushed to the surface by its mother or by an attendant female for its first breath. After the umbilical cord snaps, close to the umbilicus, the baby begins to swim. The fins and flukes are folded and flaccid at birth but harden into proper position within a few hours. The young of both porpoises and larger whales accompany the mother very closely for the first few weeks, often positioning themselves at about the midbody of the parent and taking advantage of water flow, so that they are carried along with a minimum of effort.

Baleen whales apparently breed once every two years and nurse young from seven to more than ten months. The nursing span is more protracted in toothed whales: the young bottlenose dolphin may continue to nurse up to the age of 18 months and the young sperm whale for 12 to 13 months.

Mammary glands and milk

The paired mammary glands in all whales are hidden beneath the body blubber, the nipples being buried in a longitudinal fold located lateral to the genital aperture. Milk, which is secreted into the enlarged lacteal duct, is thought to be forced into the mouth of the baby by contractions of body muscles underlying the mammary gland. In baleen whales and most porpoises, the baby thrusts its snout into the mammary slit, taking the nipple in its tongue and receiving the milk; the tongue, when applied to the palate (roof of the mouth) forms a tube for receiving the milk. The baby sperm whale, which has a bulbous snout, takes the nipple sideways in its jaw. Whale milk is very high in fat content and almost lacking in milk sugar.

Migration. Most baleen whales undergo seasonal migrations, some as long as 5,000 kilometres (3,000 miles) each way, from feeding to calving grounds. With the exception of the sperm whale, odontocete migrations seem to be much more local in character and in some species apparently do not occur at all. There is a seasonal migration of sperm whales from calving grounds near the Equator toward the higher latitudes, but only adult males go beyond the temperate zones.

Summer concentrations of planktonic (free-floating) food, such as the euphausiid crustacea, or krill, are greatest in polar latitudes; hence, baleen whales congregate in these waters to feed. Movements toward the Equator occur in the austral (southern) or boreal (northern) fall and apparently terminate within about 30° of the Equator, usually far at sea, where birth of young takes place in the warm waters. The humpback whale (*Megaptera novaeangliae*) typically frequents island shores at this time, and mothers and young may be seen especially near subtropical islands. Long migrations are apparently made wholly or nearly completely without feeding, except on the summer grounds; the animals live off body fat reserves in the blubber coat. Some baleen whales feed closer to the Equator in such areas as the Persian Gulf and the Gulf of California; migrations in these populations, therefore, are probably of lesser extent. The migratory cycles of the gray whale (*Eschrichtius robustus*) and the humpback conform closely to an annual day length cycle, suggesting that a photoperiodic effect may be involved in the regulation of reproduction.

In the smaller toothed whales, some food-related seasonal movements occur in which porpoises and whales may congregate in certain areas in response to seasonal abundances of squid or fish. Other porpoise species are sedentary, and some continuously occupy discrete portions of a shoreline. Odontocetes apparently do not undergo prolonged periods of starvation and thus must have continuously available food supplies.

Locomotion. Swimming and diving. Swimming is accomplished by vertical undulatory movements of the tail, the thrust coming from both the horizontal flukes and the posterior part of the body. Typically, a swimming cetacean describes a roughly sinusoidal (undulatory) path as it swims, coming to the surface for a breath at the

top of each curve and descending below the surface in between. Because breathing is accomplished through the paired or single blowhole on top of the head, the animal does not have to arch its neck to breathe and thus break the smooth path of travel. Only in the sperm whale is the blowhole anterior on the left side of the snout tip. The flukes are frequently thrown free of the water before steep dives, a habit uncommon in rorquals (*Balaenoptera*).

Small cetaceans normally do not dive for periods exceed-

Vincent Serventy



Figure 41: Humpback whale surfacing, showing blowholes.

ing five to seven minutes, but the beaked whale (*Ziphius cavirostris*), the bottle-nosed whale (*Hyperoodon ampullatus*), and the sperm whale may stay below for very long periods; the sperm whale may routinely dive for an hour below the surface, and specimens that have become tangled in submarine cables indicate that it may reach depths of at least 1,100 metres (about 3,600 feet).

Length of dives

Speed. It has been reported that porpoises travel at speeds of at least 38 kilometres per hour (21 knots, or 24 miles per hour) and whales at 56 kilometres per hour (35 miles per hour). Tests on trained animals suggest that cetaceans possess unusually low drag because the boundary layers (water near the skin) remain partially laminar (smooth) at high speeds and do not become predominantly turbulent, as is the case around moving rigid bodies. Body shape is believed to play an important part in the low drag experienced by swimming whales.

Leaping and wave riding. Many whales and porpoises leap from the water from time to time, and some perform acrobatic manoeuvres. Sperm whales reportedly can jump clear of the water; humpback whales may mate as they leap in tandem from the water; and many porpoise species leap during fast swimming. Some species trained in captivity have vaulted as high as six metres (20 feet) above the water surface. One species group of the genus *Stenella* spins rapidly along its longitudinal axis as it leaps. Some of these leaps are thought to aid in shaking off suckerfish (*Remora*), but use in social contexts is also evident. The larger rorquals seem never to leap; the gray whale spy hops, or pitchpoles, often rising to about the level of its flippers and hanging for a moment, apparently viewing the aerial surroundings; or it may leap three-quarters of its length and fall back on its side.

Many species of smaller whales and porpoises come to the bows of vessels and may ride for many minutes without beating their tails, obtaining a thrust from the pressure wave in front of the ship's bow wave. Porpoises have also been observed riding the "bow wave" produced by a large mysticete when the latter was at the surface.

Sound production and communication. Acoustic behaviour, which is highly developed in all cetaceans, involves passive listening, social signalling, and echolocation. Although the acoustic sense and passive listening are developed in all cetaceans, only odontocetes have been shown to echolocate—*i.e.*, to use the echoes of their own

signals for discrimination and navigation. Two major sorts of sounds are produced by cetaceans: low-pitched signals predominantly in the human audible range, such as barks, whistles, screams, and moans; and brief clicks of high-intensity sound, some as high as 200 kilohertz, or about 13 times the upper frequency of humans. The former sounds seem to serve mostly in social communication; the latter serve largely in discrimination and navigation. Click production in porpoises allows highly refined discrimination; in experiments, porpoises can discriminate between objects of slightly differing sizes, detect the presence of fine wires, determine the locations of tiny objects, and discriminate among different kinds of metals. Clicks are projected from the porpoise forehead and upper jaw in a highly structured beam; the highest frequencies are strongest directly ahead of the animal. Low-frequency sounds are more omnidirectional, and their source in dolphins is thought to be in structures within the left nasal passage.

Types of communicative signals

Communicative signals are used in a variety of contexts. Humpback whales produce a variety of moans and calls in their breeding areas, porpoises emit distress cries, and specific signals are produced during birth. The finback whale produces a monotonously repeated, very low-pitched, narrow band sound (about 20 Hertz) that is thought to allow long distance communication between whales. Captive porpoises have mimicked human voice signals and artificial sounds, and a variety of sounds have been shaped by training techniques. No evidence exists, however, for a humanoid language in any cetacean.

Nonvocal signalling of emotional states, which occurs in porpoises and small whales, has not been studied in the larger whales. Such signals include showing the whites of eyes, which is a sign of rage, and various jaw claps and body postures, which carry social meaning.

FORM AND FUNCTION

General features. The general body form is fusiform, or spindle-shaped, with the head end variously modified into a more or less attenuated beak, rounded, bluff, or flattened. The tail always ends in a horizontal blade (the flukes), through which vertical movements of the tail produce forward thrust. The forelimbs are paddle- or sickle-shaped flippers that function in balance and steering. A dorsal fin is usually present. The flippers have all the recognizable elements of the mammalian forelimb skeleton. No external trace of the hindlimbs remains, but a greatly reduced pelvis is represented by a pair of slender, irregularly curved bones, remote from the backbone and embedded in the flesh in the vicinity of the reproductive opening. In some of the larger whales, bony or cartilaginous remnants of the hindlimb skeleton still persist as attachments to the pelvis.

The head varies in size from about one-tenth of the total length in some toothed whales, such as dolphins, to nearly one half in the sperm whale. No ear pinnae are

to be distinguished adjacent to the inconspicuous opening of the external ear tube. Eyes and eyelids are present and usually functional, but vision is said to be reduced in some river dolphins.

The facial part of the skull is elongated into rostrum, or beak (Figure 42). The maxillary bones are spread laterally at their posterior ends and either override the frontal bones (in the toothed whales) or extend under the orbital processes of the frontal bones (in baleen whales). The telescoping of the skull and the laminated (layered) bony structure, together with the high, short, broad braincase, result in a skull form that is far removed from the usual mammalian pattern.

Feeding adaptations. The digestive tract and the dentition show several special features associated with feeding habits. Food enters the large gape and passes unchewed through the esophagus and into a specialized multichambered stomach. Specialization of the stomach is related to the simplification of dentition, which requires that the food, swallowed whole, undergo prolonged digestion. In the modern carnivorous cetaceans (the odontocetes), the teeth are simple, uniform in shape, peglike, and single rooted, frequently exceeding the primitive mammalian number (44). When the mouth is closed, the top teeth interdigitate with those of the lower jaw, rather than occluding, as do those of most land mammals. The teeth are used for seizing prey, which is not chewed before being swallowed. Among the dolphins, many species have reduced dentition. Risso's dolphin (*Grampus griseus*), which feeds predominantly on cuttlefish, has the teeth reduced to from two to six on each side of the lower jaw and none in the upper jaw. This reduction of teeth in association with a cuttlefish diet is even more evident in the bottle-nosed and beaked whales (Hyperoodontidae), whose slender jaws bear one pair (occasionally two pairs) of functional teeth in the mandible and none in the upper jaw, except in the Tasmanian beaked whale (*Tasmacetus shepherdi*), which has teeth in both jaws, totalling up to 90. The sperm whale, in which functional teeth are restricted to the lower jaw, eats large squid extensively, but fish are also taken.

The baleen whales are plankton feeders, including in their diet small crustaceans such as krill (euphausians) and copepods (*Calanus*), and fishes such as herring, sardines, and capelin, which they strain from large amounts of water with the baleen plates. These are anchored to the roof of the mouth, somewhat spaced apart at right angles to the long axis of the head. Each plate is triangular, with its smallest side inserted in the jaw. The outer (labial) edge of the plate is smooth; the inner (lingual) side is frayed out into a fringe of bristles that collectively form a matted sieve or strainer. The strainer selectively removes the plankton from the water expelled from the whale's mouth. The plates are longest (to nearly four metres, or 13 feet) in right whales and shorter, broader, and stiffer in orquals. Each half of the total baleen complement may comprise upward of 300 plates. Each plate is formed of parallel tubes of horn (keratin), with compacting horn between, which are enveloped, except at their fringed ends, in a covering layer. Between adjacent plates is a pad of softer horn. Baleen does not differ essentially from hair in chemical composition and, like the latter, has its origin in the skin.

The senses. *Smell.* The sense of smell is reduced or lacking in cetaceans. A reduced olfactory organ is, however, present in baleen whales; on the other hand, it is absent in odontocetes.

Sight. Underwater vision is excellent in porpoises (in species in which it has been tested), except for some river porpoises. The eyes of the Amazon River porpoise (*Inia geoffreyensis*) are reduced but apparently functional; those of the susu of the Ganges, Brahmaputra, and Indus rivers (*Platanista gangetica*) are rudimentary and vision is poor. Some evidence indicates that porpoise vision out of water is poor. The lens of the cetacean eye is modified for underwater vision: the sclerotic capsule of the eye is greatly thickened, and a variety of adaptations relating to pressure change and vision in water are present. There is no definite evidence for colour vision in cetaceans.

Hearing. Hearing is a major sense in all cetaceans.

Dentition

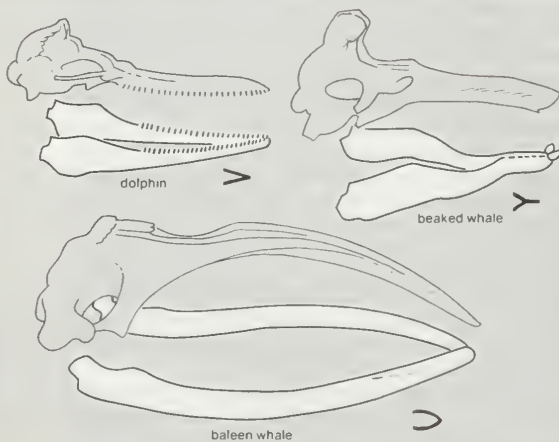


Figure 42: Three cetacean skulls. The small figures represent the general shape of the jaw as viewed from above.

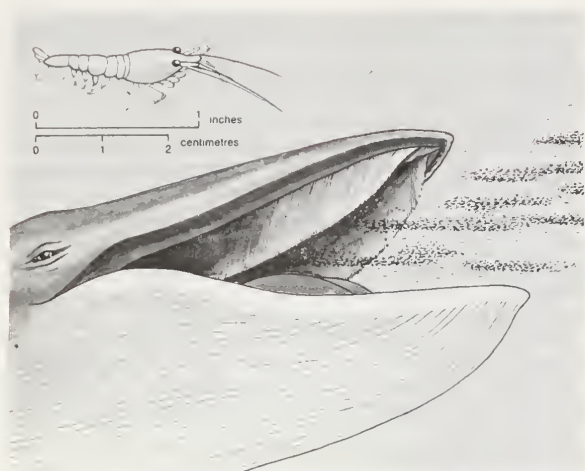


Figure 43: Fin whale swallowing krill, one of which is shown in the inset above.

In the bottle-nosed dolphin hearing capability extends through the human range (20 Hertz to 20 kilohertz) to about ten times the human upper limit (up to about 150 kilohertz). In odontocetes, the hearing of clicks is highly directional and is apparently mediated through the lower jaws and throat, rather than through the much reduced external ear canal. Three species of Japanese porpoises are approximately 10 times as capable as a human of discriminating pairs of clicks given in rapid succession.

Little is known about the sound reception of baleen whales other than evidence from whalers who know them to be extremely sensitive to boat noise.

The ear cartilages, of similar origin to those that support the pinna (external ear) in land mammals, have been sunk in the blubber over the external ear. In baleen whales, the external ear canal becomes a solid cord; it opens again near the eardrum and is filled with a solid waxy plug the laminations of which have been used for age determination, much in the manner of tree rings. In odontocetes, the tiny canal is generally open during its entire convoluted length. In one species, the Dall porpoise (*Phocoenoides dalli*), it is at least occasionally occluded at the surface. The middle ear ossicle bones differ widely from those found in terrestrial mammals, apparently in response to the conditions of hearing underwater. The middle ear bones and the inner ear (cochlea) are housed within a heavy, almost ivory-like, petrotympanic complex, which hangs nearly free from the skull and is surrounded by air sacs connecting with the respiratory tracts by way of the eustacian tube. These sacs are important in pressure equilibration during dives and in acoustic isolation of one ear from the other and of both ears from the animal's own sound signals. Masses of small blood vessels occupy much of the middle ear and are also instrumental in pressure equalization.

Taste. Marked food preferences have been noted in captive odontocete cetaceans. The presence of well-developed taste buds on the tongue point to a well-developed sense of taste.

Proprioception. The awareness of bodily position (proprioception) is as acute in cetaceans as in any other mammal. Their exquisitely refined capabilities in manoeuvring in the water provide other evidence of coordinative sense.

(K.S.N.)

PALEONTOLOGY AND EVOLUTION

It is generally agreed that whales originated in the early Paleocene or Upper Cretaceous (about 70,000,000 years ago) from some group of terrestrial meat-eating mammals. Some authorities believe that a likely parent group was the Mesonychidae, medium-sized, large-headed mammals of the order Creodonta, believed to be near the base stock of both the carnivores (Carnivora) and the ungulates (several orders of hoofed mammals). It has also been suggested that the whales evolved from a Paleocene member of the order Insectivora resembling the Eocene *Pantolestes*, a semi-aquatic carnivorous animal about the size of the modern otter (*Lutra lutra*). An American paleontologist, L. Van

Valen, has pointed out similarities between the mesonychids (as well as the closely related arctocyonids) and the earliest archaeocete whales, the protocetids. All that can be stated with certainty is that the ancestors of the protocetids must have evolved through an amphibious stage.

The earliest known whale fossil is a fragmentary shoulderblade from the lower Eocene of England (about 50,000,000 years ago). Remains of archaeocete whales of the family Protocetidae are found in middle Eocene formations. They were relatively small animals about the size of the modern porpoises (two to three metres long). The hindlimbs had already been lost, and the nostrils were midway between the terminal position found in most terrestrial vertebrates and the top of the head, where they are found in most advanced whales. The snout was elongated, but the skull bones were not telescoped as in modern forms; the teeth retained the differentiated (heterodont) condition found in terrestrial carnivores, the front teeth being tapered pegs and the cheek teeth serrated. At the time of their greatest development, in the late Eocene, some archaeocetes were large, the most striking being *Basilosaurus*, a slender animal that attained a length of some 20 metres (nearly 70 feet). One line of archaeocetes, the dorudontids, persisted into the early Miocene (about 20,000,000 years ago).

Paleontologists disagree as to whether the modern toothed whales (odontocetes) and baleen whales (mysticetes) arose from an early archaeocete relative or separately from different terrestrial ancestors. Recognizably odontocete remains are known from the upper Eocene, and even earlier odontocetes must have been contemporaries of the early archaeocetes. One early form, *Agorophius*, already showed some telescoping of the skull and the nostrils were above the eyes. The agorophiids appear to have had a relatively short history but were probably ancestral to the porpoise-like squalodonts, which appeared in the late Oligocene (about 28,000,000 years ago) and by the middle Miocene were the dominant toothed whales of the world's oceans. By the early Miocene most of the modern odontocete families were represented by some members. Fossils of sperm whales found in lower Miocene deposits have been assigned to the modern genus *Physeter*. These are among the oldest known representatives of the physeterid line, which produced a variety of other Miocene and later forms and must have existed as a well-differentiated group during the later part of the Oligocene.

Two groups of fossil whales appear to be intermediate between the toothed and the baleen whales. In the Patriocetidae the nostrils were in an intermediate position between the snout and the forehead; the telescoping of the skull (in the symmetrical manner seen in later mysticetes) was only slightly evident. The other group, the Aetiocetidae, was placed in the Odontoceti on the basis of possession of teeth, but Van Valen has pointed out that *Aetiocetus* is mysticete in other respects and should be considered closer to the mysticete evolutionary stem.

The earliest fossils definitely assignable to the Mysticeti are the middle Oligocene cetotheres, a group that became quite abundant in the subsequent Miocene. Of the three modern families of baleen whales, two (Balaenidae and Balaenopteridae) are known from scanty Miocene and abundant Pliocene fossils; the third (Eschrichtiidae) is known only from a few Pleistocene and post-Pleistocene fossils representing the sole modern genus, *Eschrichtius*.

CLASSIFICATION

Distinguishing taxonomic features. The three suborders recognized below share the same basic body plan but differ in the degree of specialization attained. The suborders and families are separated primarily on the basis of the following characteristics: tooth structure, number, and degree of differentiation; modifications of the skull, especially the position of the nostrils, degree of telescoping of the whole skull, and extent of joining of the two halves of the jaw (rami); and degree of modification of the pelvic girdle (in the archaeocetes only).

Annotated classification. The classification presented below is based on research by cetologists F.C. Fraser, R. Kellogg, and a number of other modern authorities.

Structure
of the
ear

Inter-
mediate
fossil
forms

Groups marked with a dagger (†) are extinct and known only from fossils.

ORDER CETACEA

Aquatic mammals with forelimbs modified into flippers and hindlimbs lacking; pelvic girdle vestigial and not attached to vertebral column; tail laterally flattened and extended into horizontal flukes, supported by cartilage. External nares at top of head (except in sperm whales and archaeocetes); cranium telescoped variably, with elongated rostrum. Three suborders, with about 10 fossil and 9 Recent families.

†Suborder Archaeoceti (archaeocetes or Zeuglodonts)

Fossil only; lower Eocene to middle Miocene. Anterior and posterior teeth differentiated; total teeth not exceeding 44, the basic number in terrestrial placental mammals; nostrils positioned toward the top of the head, to varying degrees. More than 12 genera described.

†Family Protocetidae

Lower to upper Eocene; Europe, Africa, and possibly North America.

†Family Dorudontidae

Upper Eocene to lower Miocene; Europe, Africa, North America.

†Family Basilosauridae

Upper Eocene to lower Oligocene; Europe, Africa, North America.

Suborder Odontoceti (toothed whales)

Upper Eocene to Recent; worldwide marine, a few in fresh water. Teeth relatively uniform (homodont) and numerous (up to 300), occasionally reduced. Skull asymmetrical; nostrils on top of head (except in Physeteridae) with single external opening (blowhole). Several pairs of ribs joining sternum, which is fused into a single bone only in the adult.

†Family Agorophiidae

Fossil only; upper Eocene of North America. Two genera.

†Family Squalodontidae

Fossil only; upper Oligocene to lower Pliocene; Europe, North and South America, Australia and New Zealand. More than 12 genera described.

Family Platanistidae (river dolphins)

Lower Miocene to Recent. Snout long and slender; neck vertebrae all free; sight reduced. Flippers short and broad, sternum well developed. About 10 fossil and 4 monotypic Recent genera, the latter respectively in the Ganges, Indus, and Bramaputra rivers of India; La Plata of Argentina; Amazon and Orinoco of South America; and Tung-Ting Lake and adjacent Yangtze River, China. Adult length 1.5–2.5 metres (5–8 feet).

Family Hyperoodontidae (beaked whales)

Previously called Ziphiidae. Lower Miocene to Recent. Functional teeth greatly reduced in number, to 1 or 2 pairs in lower jaw (except for total of 90 teeth in *Tasmacetus*); vestigial (non-functional) teeth not uncommon. Pair of longitudinal grooves in throat region. Dorsal fin slightly more than one-third of way back along body. Tail usually not notched in middle. Flippers small. Worldwide, marine; length to about 12 metres (40 feet). About 14 fossil and 5 Recent genera; 15 Recent species.

Family Physeteridae (sperm whales)

Lower Miocene to Recent. Head proportionately large, with bulbous, squared snout; mouth narrow and ventral; lower teeth totalling 40–52, upper teeth vestigial, smaller, varying in number. Blowhole a single, asymmetrical opening on the anterior left tip of snout. Twenty fossil and 2 Recent genera; length to 19 metres (62 feet; *Physeter*) and to 3 metres (10 feet; *Kogia*) in Recent genera. Tropical and temperate oceans, and (adult males only) polar waters.

Family Monodontidae (narwhal and beluga)

Pleistocene to Recent. Dorsal fin lacking. Neck vertebrae free. Flippers broad, rounded at tips. Teeth reduced to 8 or 10 in *Delphinapterus* (beluga); all vestigial in *Monodon* (narwhal) except for 1 left tooth of male, which grows into a long, straight tusk, extending in front of animal. Two genera; length to about 5 metres (16 feet); Arctic Ocean.

Family Phocoenidae (porpoises)

Upper Miocene to Recent. Head with rostrum (forehead extending to snout); teeth with expanded spade-shaped crowns. Dorsal fin triangular (when present); tail notched posteriorly. Three fossil and 3 Recent genera. Adult length (in Recent species) to about 2 metres (6.6 feet); virtually worldwide; marine, 1 species ranging up Yangtze River.

Family Delphinidae (dolphins and killer whales)

Lower Miocene to Recent. Characteristics variable. Dorsal fin

and well-defined beak present or absent. None with expanded tooth crowns. About 35 fossil and 12 Recent genera. Length 1.7 to about 9 metres (5.5 to 30 feet), worldwide, marine; 1 species in salt and fresh water (Irrawaddy River).

Family Stenidae (long-snouted dolphins)

Lower Miocene to Recent. Like Delphinidae but differing in structure of air sinus system and in shape of forehead. Three genera. Widespread in tropical oceans and rivers.

†Families Eurhinodelphidae, Hemisyntrachelidae, and Acrodelphidae

Three families of Miocene toothed whales; about 12 included genera, from all continents.

Suborder Mysticeti (baleen, or whalebone, whales)

Teeth lacking in adult; numerous but vestigial in embryo; baleen present. Lower jaw large, mandibles bowed outward and loosely united in front. Skull symmetrical; nasal bones relatively well developed; blowholes paired slits, slightly divergent posteriorly.

†Family Aetiocetidae

Upper Oligocene. Toothed but with symmetrical skull and other typical mysticete features. One genus (possibly 2); North America.

†Family Patriocetidae

Upper Eocene to upper Oligocene; Europe and North America; 4 genera.

†Family Cetotheriidae (cetotheres)

Middle Oligocene to lower Pliocene; North and South America, and Europe; about 30 genera.

Family Eschrichtiidae (gray whale)

Formerly Rhachianectidae. Head less than one-quarter of total length; neck vertebrae not fused. Dorsal fin lacking. Flippers with 4 digits. Pelvis large. One Recent species; length to about 15 metres (49 feet); found along both coasts of the North Pacific Ocean. Fossils from Pleistocene and Recent, from Atlantic Ocean.

Family Balaenidae (right whales)

Lower Pliocene to present. Skull much arched; neck vertebrae fused; baleen long, narrow, flexible strips. Ventral grooves on throat lacking. Four fossil and 2 Recent genera (4 species); length 6–20 metres (20–66 feet); polar and subpolar oceans.

Family Balaenopteridae (rorquals and humpback)

Upper Miocene to present. Skull broader and less arched than in Balaenidae; baleen plates shorter, broader, less flexible, neck vertebrae not fused. Dorsal fin present; flippers narrow. Conspicuous longitudinal grooves present on throat. Eight fossil and 2 Recent genera (8 species); length 10–33.6 metres (33–110 feet); blue whale (*Balaenoptera musculus*) is the largest animal that has ever lived. Family occurs in all oceans and, at various times, from Equator to poles; migratory.

Critical appraisal. Authorities disagree on whether the whales should be treated as one order, Cetacea (as above), two, or three (as in the classification used in the section *The class Mammalia*). The determining factor is the degree of shared ancestry, the evaluation of which remains controversial. There is a possibility that certain protocetids could have given rise to both the modern groups of whales, and it is because of this that some authorities prefer the use of a single order. However, the absence of intermediate fossils linking the baleen whales with the toothed forms supports the use of separate orders. Resolution of this problem must await the discovery of relevant fossil material.

Below the ordinal level, there is further disagreement at many points. Although there is no doubt that any recent cetacean is odontocete or mysticete, the relationships of many genera are in doubt. The long-snouted dolphins (Stenidae) are of uncertain affinities but are included in the Delphinidae by many authorities, as are the porpoises (Phocoenidae). At the species level, there is uncertainty about the specific or subspecific status of many populations. (Ed.)

Proboscidea (elephants)

The order Proboscidea comprises three suborders and about 300 species of terrestrial mammals. All but two species, the Asiatic, or Asian, elephant (*Elephas maximus*) and the African elephant (*Loxodonta africana*), are extinct. The elephants are the largest surviving land ani-

mals and, among the mammals, are exceeded only by the whales in size.

The Proboscidea are characterized by columnar limbs, bulky bodies, and elongated snouts. In recent forms, testes are internal. The snout is a long boneless proboscis, or trunk; it is a combination of the upper lips, palate, and nostrils. Some of the incisor teeth develop into tusks. One extinct suborder (Deinotherioidea) lost the upper tusks; certain others have lost the lower ones and evolved upper tusks of dentine from which the enamel has partially or completely disappeared. The canine teeth were generally repressed in all groups, and the cheek teeth developed rows of blunt cones or ridges. In later forms, the temporary teeth were replaced by permanent ones, which are pushed by an escalator-like movement along a horizontal plane, so that the front teeth were replaced by teeth moving forward from the rear. The skull, which originally was elongated, became shorter, higher, and bulkier in later forms. The back of the eye orbit remained open instead of forming a complete bony ring, and the nasal opening in all Proboscidea is at a higher horizontal plane than the eye sockets. The neck shortened as the animals evolved larger, higher bodies and an elongated trunk that also functions as a hand. The skull has enlarged out of proportion to the brain in order to serve as an anchor for the trunk and to support the heavy dentition. This order occurs in all the continents except Australia. Fossils of proboscideans provide valuable information about early humans who were their contemporaries.

GENERAL FEATURES

Size range and distribution. In Europe, the landmass that broke up to form the islands of the Mediterranean Sea harboured proboscideans. Three fossil species have been found in Malta; one had had a height of 2.1 metres, or 6.9 feet (*Palaeoloxodon mindriensis*), another a height of 1.5 metres, or 4.9 feet (*P. melitensis*), and the third was less than one metre, or about three feet (*P. falconeri*). *P. creticus* of Crete was 1.5 metres, and *P. cypriotis* of Cyprus was 0.9 metre (three feet) in height. In North America, a small *Mammuthus* isolated on Santa Rosa Island, off the coast of southern California, was probably derived from *Mammuthus meridionalis*, a species that stood 4.2 metres (13.8 feet) at the shoulders.

Early human records of elephants

Elephas maximus asurus lived in Iran and Syria. Early drawings of the animal and fragmentary skeletal remains indicate that it was the largest subspecies of the Asian elephant. The war elephants employed by Pyrrhus in 255 BC and engraved upon Roman seals show animals of unusual size. "Sarus," which signified "the Syrian," was the outstanding animal in the elephant battle squadron of the Carthaginian general Hannibal. In 1500 BC elephants (*Elephas maximus rubridens*) existed in China as far north as Anyang, in northern Honan Province. Writings from the 14th century state that elephants were still to be found in Kwangsi Province.

Man as well as other environmental factors exterminated the woolly mammoth (*Mammuthus primigenius*) and the imperial mammoth (*M. imperator*) about 10,000 years ago. Several races of the living species of Asian and African elephants also died out by about 1500 BC. The small North African race became extinct by the 2nd century AD, and some of the American mastodons, such as *Cuvieronius postremus* of South America, died out as recently as the 4th century AD. The large African bush elephants (*Loxodonta africana*) were exterminated from the Transvaal in South Africa early in the 20th century, but they still occur over much of the continent south of the Sahara Desert. A smaller elephant inhabits the forests of western equatorial Africa, particularly in the Congo region. It is considered by some to be a subspecies (*Loxodonta africana cyclotis*) of the African elephant; others believe it to be represented by several subspecies; still others consider it to be a separate species (*L. cyclotis*).

In Asia, elephants have been exterminated from Iran, Iraq, Afghanistan, the northwestern part of India, and from much of the Malay Peninsula, Java, and the greater part of Borneo and Sri Lanka (formerly Ceylon). Isolated colonies remain in forest areas of Mysore in the peninsular

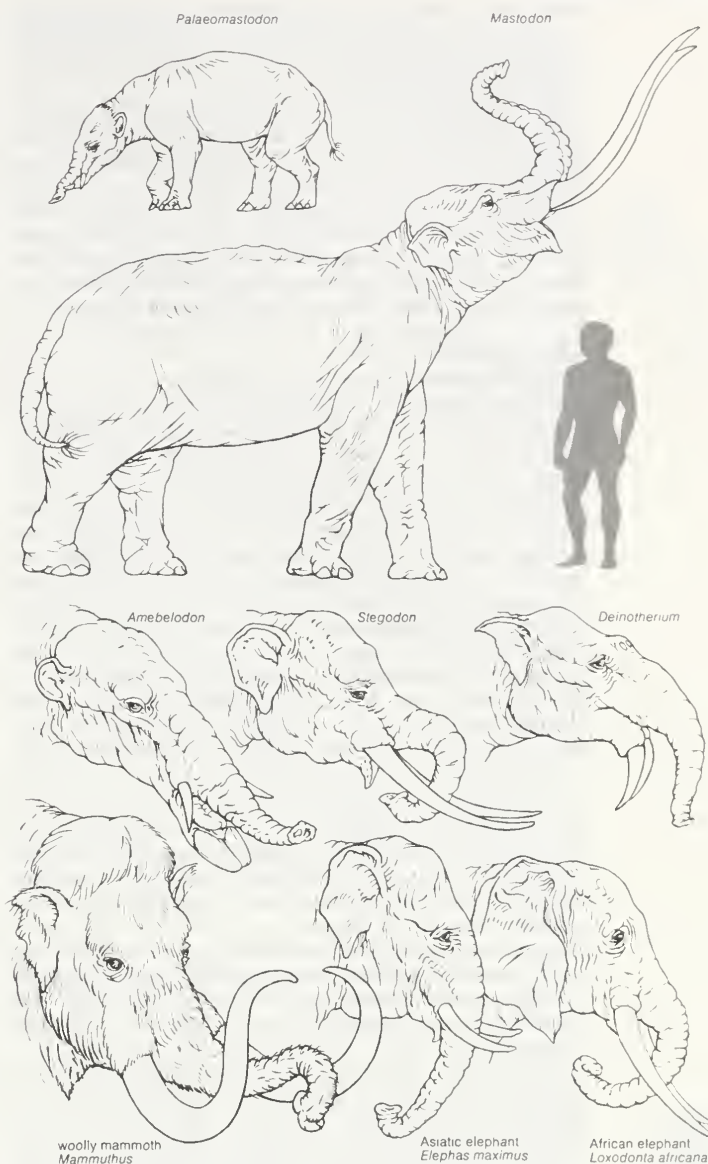


Figure 44: Representative fossil and living proboscideans. Drawing by Christian D. Olsen

part of India, in Assam, Nepal, Burma, Malaya, southern China, Sri Lanka, Sumatra, Borneo, and other islands of the East Indies.

Importance to man. Elephants constitute the chief source of commercial ivory. Because of the continuing demand for this commodity, the animals are in danger of extermination. Elephant "pearls" consist of concentric layers of ivory deposited over a foreign object that has been intruded into the soft ivory at the base of the growing tusk. From early times elephants have been used as beasts of burden in India and Burma. Since they do not breed freely in captivity, new stock for domestication is often captured from wild herds. One method is to drive them through a funnel-shaped stockade into a small enclosure; trained tame elephants help to subdue, noose, and train young captives for service. The training and handling of an elephant are usually entrusted to one man, called the mahout in India and the *oozie* in Burma; an elephant and its keeper frequently became inseparable companions.

Capture of elephants

Trained elephants carry humans in a howdah, or miniature hunting lodge, on their backs. They are used to move timber or other heavy materials.

In modern times, African elephants have been trained for labour only since late in the 19th century; however, they were used by the Carthaginians in wars with the Romans. Both species were tamed at that time. Elephants were used as executioners, in amphitheatres, and for military pageants. They are still used in exhibitions at circuses, carnivals, and zoological parks. In warfare elephants have

been used to drag heavy equipment, especially through mud and up steep slopes. As late as World War II, they were of value in military movements over the mud-clogged roads of Southeast Asia.

The association between man and elephants goes far back into mythology, and a rich folklore has developed. Bracelets of hairs from the tail of a freshly-killed or living elephant are prized as good-luck charms.

Attempts to locate the legendary "graveyards" to which old elephants allegedly resort when near death have been, for the most part, unsuccessful. The groups of buried elephant remains that have been found probably represent sites where elephants drowned in bogs or quicksands.

NATURAL HISTORY

Reproduction and life cycle. Tuskers, or tusked bulls, occasionally fight brief, savage duels that may end in death for the defeated animal. A duel between tuskless elephants may last for days, with occasional periods of rest. The female selected from the herd by the winner often makes an apparent attempt to escape from him. After a brief preliminary courting, the male mounts the female from behind, leaning over her back and either gripping her body or resting his forefeet upon her pelvis, and assumes a standing posture. Copulation lasts for about 20 seconds, with very little movement or noise. Mating continues promiscuously for about two days, after which the most powerful bull drives off the others and remains with the cow for about three weeks. The period of gestation varies from 20 months for a female calf to 22 months for a male. When parturition is about to occur, the herd surrounds the cow, who assumes a squatting position while giving birth.

In regions where large carnivores, such as tigers, prey upon newborn elephants, the cow seeks a female associate. The mother and the other elephants in the herd blow dust upon the moist, newborn calf to dry it. Two hours after birth, the baby is able to stand and is suckled. The mother and calf then join the herd.

Tame cows begin breeding at the age of eight or nine years; tame bulls begin when about 11 or 12 years of age. The interval between the birth of successive calves is about four years. In captivity cows are known to continue bearing calves until 60 or 70 years of age.

The newborn elephant is about one metre (three feet) high and weighs about 90 kilograms (200 pounds). It is covered with yellow and brown hair. After a few months the hair on some parts of the body is as long as that in the extinct mammoth. In *E. maximus* there is also a pinkish patch around each eye, and when the calf is about five months old, a faint whitish patch develops on each cheekbone. As this patch spreads, similar patches appear upon the trunk and ears. In the more easterly races of Malaya and Sumatra there are only a few gray spots. The hoof nails, which are dark at first, gradually become lighter. When the elephant is about eight years old, a thick oily secretion known as musth, or must, begins to flow from a gland in the temporal, or temple, region. It occurs in both sexes, but is more active in males. The secretion increases each year until it drips into the elephant's mouth. Some authorities believe that the function of the secretion is to inhibit feeding; others believe it has some effect on sexual activity.

An elephant is not fully grown until it is about 25 years old. In the wild, the average life-span is about 60 years, but under optimum conditions an elephant may live for 80 years.

Behaviour. The organization of an elephant herd is often according to age and sex. In *Elephas*, although herds of 10 or 15 females and their young may appear to be under the leadership of a large female, and their organization matriarchal, young adult males are always in the vicinity, as is the real leader of the entire group, a large male. The leader may be accompanied by one young adult male who acts as a scout, warning the leader of danger. The herd also has a system of scouts, and, before emerging into an open area, one of the scouts usually explores it. If no danger is apparent, he signals by trumpeting to the herd to advance. Individuals often serve as guards while the rest of the herd feeds or bathes.

The herd is held together both by blood relationship and by a strong sense of companionship. If an individual is injured, three or four others surround it, shielding it from danger, supporting it, and helping it to move away. A calf that has lost its mother is adopted by the other cows in the herd even if they have their own calves to raise.

Ecology. Elephants clear paths through forests that are too dense for other animals. Many modern roads in elephant-inhabited countries originated in this manner. Elephants browse to a height of about five metres (16 feet), thereby increasing the amount of sunlight available for shrubs. Their uprooting of grass and roots aerates the soil and stimulates the growth of plants that replace the ones devoured. Mud wallows frequented by elephants are fertilized by their excreta.

An elephant may destroy or discard as much vegetation as it consumes. An adult may eat between 250 to 350 kilograms (550 to 770 pounds) of solid food each day. When grazing, the animal uses its trunk or forefoot to gather grass, which is slapped against a forelimb to rid it of sand. In rainy weather, when soil is more difficult to shake off, the animal browses. Asian elephants break off branches; African elephants are more likely to push over trees. The wood apple (*Feronia elephantorum*) is a favourite food of the Asian elephant. The animal also eats wild rice that grows in forest lakes and various other aquatic plants. The African bush elephant eats the fruit of various palm trees. The spongy wood of the baobab tree provides some water during periods of drought.

FORM AND FUNCTION

Extant forms. The adult elephant has a tuft of hair at the tip of the tail and sometimes a patch of hair on the head. The limbs are adapted for bearing great weight: the legs are straight and pillar-like, and the bones of the joints are flat at the articular surfaces. Each toe has a heavy hoof nail; the weight is borne on thick pads behind the toes. The nose and upper lip are extended into a long trunk, which contains the nasal passages and has nostrils at the tip. Water for drinking is sucked into the trunk and then discharged directly into the mouth. The trunk is used for placing food into the mouth, for spraying and dusting the body, for lifting or moving heavy objects, and even for throwing objects at man.

The upper second incisors are typically developed into ivory tusks, the longest and heaviest teeth of any living animal. Canine teeth are absent. The complex molars are of the high-crowned type, with transverse rows of enamel ridges on the grinding surface, which is often traversed by a longitudinal median groove. There is normally only one complete functional tooth and half of a second one at a time on each side of each jaw. These are replaced horizontally from the rear as they wear away.

The Asian elephant (*Elephas maximus*) and its races are distinguished from the African elephants by being somewhat smaller and by having relatively small ears with the upper edge curled forward. The head is more domed, is structurally more complex, and has a greater development of diploe, or bony cell cavities. *E. maximus* also has high-crowned teeth and a relatively smooth trunk with a single fingerlike projection at the tip. Adult bulls weigh up to 5,450 kilograms (six tons) and commonly stand three metres (10 feet) at the withers. Only the males develop tusks, which average 1.5 metres (4.9 feet) in length, the pair weighing about 32 kilograms (70 pounds). These develop flat, longitudinal planes of wear near the tip. Tusks 2.7 metres (8.9 feet) long and 68 kilograms (150 pounds) in weight have been recorded. About 90 percent of the males of the Ceylonese race lack tusks; Sumatran elephants are of slighter build and have longer trunks.

In contrast to the Asian species, the African bush elephant is generally larger. This and the forest form have extremely large ears (one metre in breadth), with the upper edge curled backward, and a roughened, heavily ringed trunk with two projections at the tip. The molars are of coarser construction, have fewer ridges, are less crenulated (scalloped), and consist of a thicker surface of enamel over thick plates of dentine. The top of the head is not domed, and the forehead is more convex.

Feeding habits

Appearance of elephant calf

The tusk tips are usually conical, the legs are longer, and the eyes are relatively larger than those of the Asian elephant.

The average height of adult bull bush elephants is 3.3 metres (10.8 feet) at the shoulder, and the average weight is six tons. Cows are about 0.6 metre (two feet) shorter. Both sexes possess tusks, which average about 1.8 metres (5.9 feet) long, the pair weighing 36 to 55 kilograms (79 to 121 pounds). A pair of tusks in the British Museum weighs about 133 kilograms (293 pounds); the larger one of the two is 3.5 metres (11.5 feet) long, with a basal circumference of 46 centimetres (18 inches). The largest elephant on record, a bull bush elephant killed in the Cuando river district of southeastern Angola in 1955, is on display at the Smithsonian Institution in Washington, D.C.; it probably weighed 8,200 kilograms (nine tons) when alive and stood four metres (12 feet) at the shoulder. Adult forest elephants are about 2.1 metres (6.9 feet) tall and weigh 1,225 kilograms (2,700 pounds), with slender tusks that are often more than three metres long. The ears are relatively small and smoothly rounded at the margins.

Albinos

Albino, or "white," elephants occasionally occur, especially in Thailand and Burma, where they are regarded by some as semisacred.

The mammoth. The genus *Mammuthus* contains some of the largest members of the family Elephantidae. Certain ones reached a height of over 4.2 metres (13.8 feet) at the shoulders. One was *M. meridionalis* of Asia and Europe, and the other was *M. imperator*, which entered North America during the upper Pleistocene.

The genus also contains one of the most specialized members of the family, the woolly mammoth, *M. primigenius*, which probably became extinct about 10,000 years ago. It inhabited the sub-Arctic area of Asia and Europe and eventually entered North America over the Bering Strait; it travelled southward across western North America almost to Wyoming, then spread eastward toward Lake Michigan. Its height of 3.3 metres (10.8 feet) at the shoulders equalled that of a large *Elephas maximus*, but the body was relatively shorter, and the hindquarters sloped downward. Its skull was compressed from front to back. As an adaptation to its cold environment, the woolly mammoth evolved small ears, a short goatlike tail, and a coat of dense, furry, short hair overlain by longer, bristly hair. It also had a humplike reserve of fat upon the top of

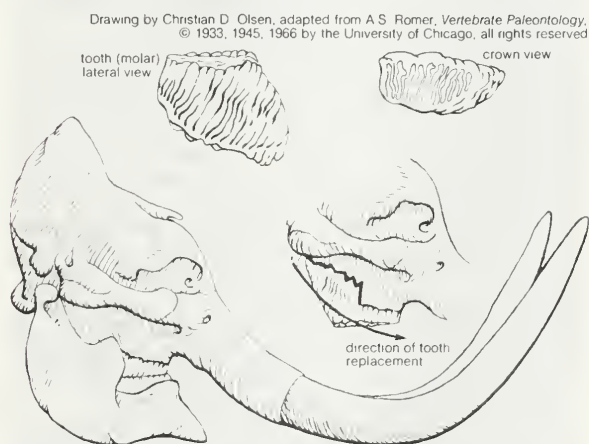


Figure 45: Skull and tooth movement in the woolly mammoth (*Mammuthus primigenius*).

the head and on the shoulders. A subcutaneous layer of fat about eight centimetres (3.1 inches) thick covered the body. The molars had as many as 27 lamellae, or plates. The tusks of the male were about 4.8 metres (15.7 feet) long. Their almost circular curvature and great size suggest that they may have functioned as shovels for exposing vegetation buried under snow. The mammoth is one of the few extinct proboscideans in which the carcass has often been completely preserved. Among the best known are two from Siberia—one discovered near the mouth of the Lena River in 1804; the other in the bank of the Beresovka River in 1900.

EVOLUTION AND CLASSIFICATION

Historical development and paleontology. The order Proboscidea has evolved from unknown ancestors that were not much larger than pigs. They flourished during the Paleocene Epoch (65,000,000 to 54,000,000 years ago). During the course of evolution, the lower jaw elongated beyond the upper, and the tusks projected well beyond the upper ones. At this stage the nose, palate, and upper lips developed into an elongated fleshy cover to the projecting lower jaw. It is probable that the nostrils opened well above the extremity of this flap and were near the eyes. A longer lower jaw proved less useful than a shorter one, so the upper flap was converted into a multipurpose tubular proboscis. Because the nostrils shifted to the tip of the proboscis, the animal was able to breathe while submerged in water. When so submerged, the animal had to rely on its sense of smell more than sight to detect approaching predators.

Origin
of the
proboscis

The suborder Deinotherioidea, consisting of one genus, is an early branch of the main proboscidean stock of the Eocene Epoch (54,000,000 to 38,000,000 years ago). They lost their upper tusks and developed a downward-hooked, tusk-tipped mandible. Numerous species of the suborder occurred in Asia, Europe, and Africa and persisted into Pleistocene times (2,500,000 to 10,000 years ago). The largest was the European *Deinotherium gigantissimum*, which reached a height of 3.8 metres (12.5 feet) at the shoulders and existed during the Pliocene Epoch (7,000,000 to 2,500,000 years ago).

In the suborder Mastodontoidea, the family Gomphotheriidae comprises 15 genera, including the earliest members of the order, *Phiomia* and *Palaeomastodon*. The former were the size of donkeys, but the latter were as large as a modern Asian cow elephant. In this family the skull and neck are elongated, and the teeth low crowned. The second incisors are enlarged; the upper ones are compressed and vertical, and they retain a band of enamel. In the later evolved genera, the lower pair are bent forward, depressed, and expanded into shovellike structures that do not meet the upper tusks. The canines are absent. Among this family are *Cordillerion* of North America and *Cuvieronius* of South America. The latter became extinct as recently as AD 200 to 400.

Phiomia, which occurred in Egypt and India, was an archaic shovel-tusked form with an elongated neck. It flourished during the Oligocene (38,000,000 to 26,000,000 years ago). The mandible and its tusks became more shovellike in *Amebelodon* and *Platybelodon* of the Miocene (26,000,000 to 7,000,000 years ago).

In *Palaeomastodon* the skull shortened, and the tusks assumed a cylindrical shape. The genus occurs in middle Oligocene deposits of Egypt. In some later genera, such as *Anancus* and *Stegomastodon*, the jaws shortened and the lower tusks disappeared. In some shovel-tusked, the mandible remained enlarged and continued to function as a shovel for digging plant bulbs. It was probably protected by a horny pad.

Mastodontidae contains the single genus *Mastodon*. It lacked lower tusks. The tusks were about two metres (seven feet) long. Some species attained a height of about three metres (10 feet), and were covered with hair. Parts of carcasses of the recently extinct American species *Mastodon americanus* have been discovered in peat deposits and swamps.

The earliest proboscideans in Asia are *Phiomia*, *Gomphotherium*, *Platybelodon*, and *Serridentinus*, of the family Gomphotheriidae, and *Stegolophodon*, of the family Elephantidae. All occurred in the Miocene of China and some in Burma and India.

In the suborder Elephantoidea, *Stegolophodon* is intermediate between the mastodons and elephants proper. Although this genus first appeared during the Miocene of Europe, Asia, and Africa, it persisted into the upper Pleistocene of these continents.

Stegodon, which occurred during the Pliocene in Asia and during the Pleistocene in Africa and Asia, evolved from *Stegolophodon*. The skull enlarged, the jaws shortened, and the lower tusks disappeared. *Stegodon zhaolongensis* of China was a large species with the most primitive teeth

in the genus. The tusks of some species such as *Stegodon magnidens* of India were about 3.3 metres (10.8 feet) in length; in others, however, they were feeble. The molars were usually low crowned, and the depression between each pair of dental folds or plates was Y-shaped, rather than V-shaped or U-shaped as in other elephants.

The genus *Mammuthus* includes all species formerly placed under *Archidiskodon*, *Metarchidiskodon*, *Parelephas*, and *Stegoloxodon*. Its species occurred in the Pleistocene of Asia, Europe, America, and Africa. Several species had great bulk and heavy tusk development.

On various occasions during the Pleistocene Epoch, normal-sized elephants that inhabited a continent were isolated when a part of the landmass was separated into islands by the submergence of low-lying land. As these isolated colonies of elephants multiplied, their fodder supply decreased, and the animals gradually became smaller—an unsuccessful measure against extinction. This process is evident in the East Indies and Philippines, in the islands of the Mediterranean Sea, and in certain islands off the coast of southern California.

Classification. *Distinguishing taxonomic features.* The proboscideans are classified largely according to body size; shape of the skull; dentition; and the shape, size, and degree of reduction of enamel in the tusks. Groups marked with a dagger (†) are extinct and known only from fossils.

Annotated classification.

ORDER PROBOSCIDEA

Oligocene to present; North America, Eurasia, and Africa. Heavy-bodied (graviportal) animals with snout prolonged into a fleshy proboscis (trunk). Tusks developed from upper or lower incisors or both; canines absent; cheek teeth with transverse rows of blunt cones or ridges. Skull short, high; nasal openings at a higher horizontal plane than eyes. Body size medium to large; shoulder height from about 1 m (about 3 ft) to more than 4 m (13 ft). About 300 species, all extinct but two.

†Suborder Deinotherioidea

†Family Deinotheriidae

Lower Miocene to upper Pleistocene; Europe, Asia, Africa. Upper tusks lacking; lower tusks curving downward from tip of lower jaw. One genus (*Deinotherium*); many species; height to about 3 m (10 ft).

†Suborder Mastodontoida

†Family Gomphotheriidae

Lower Oligocene to Recent (AD 200–400); Europe, Asia, North and South America. Skull and neck elongated. Teeth low-crowned; succession vertical. Later genera with protruding, shovellike lower incisors, others with substantial upper tusks and no lowers; premolars with 2 transverse crests, molars with 3. About 15 genera, several dozen species; shoulder height about 1 to 3 m (3 to 10 ft).

†Family Mastodontidae (mastodons)

Lower Miocene to upper Pleistocene (possibly to early historic times); Europe, Asia, Africa, North America. Molars with rounded prominences, but no ridges; lower tusks absent, but upper incisors substantial, reaching over 2 m (6.6 ft) in length in males of some species. One genus (*Mastodon*), many species; shoulder height to at least 3 m (10 ft).

Suborder Elephantoida

Family Elephantidae (elephants and mammoths)

Lower Miocene to present; fossils from Europe, Asia, East Indies, Africa, and North America; Recent species from Africa (*Loxodonta*) and southern Asia (*Elephas*) which have probably evolved from *Stegolophodon*. Six genera, with several dozen fossil and 2 Recent species; shoulder height from 1 to about 3.5 m (3 to 11.5 ft). The epiphyses (ends of long bones in limbs) do not fuse until the last molars appear. The molars are replaced at least three times. Marrow disappears from the limb bones early in adulthood.

(P.E.P.D.)

Sirenia (dugong, manatees)

Sirenians are a group of large aquatic mammals that have become rare or extinct as a result of exploitation by man for meat and oil. The largest, Steller's sea cow (*Hydrodamalis gigas*), which reached lengths of about eight metres (about 26 feet), was eaten out of existence by hungry seal hunters within a few decades of its discovery in 1741 in the Bering Sea. The remaining forms, to which the

common name sea cow is also sometimes applied, are the dugong (*Dugong dugon*) and manatees (three species of *Trichechus*), which, if their stocks were allowed to rebuild themselves, could again become of economic importance. Being the only large aquatic herbivores, other than some turtles, they could provide meat from the vast expanses of marine and freshwater vegetation, at present quite unused by man, and so bring another marginal area into production.

Dugongs and manatees are both usually seen up to lengths of three to four metres (10 to 13 feet), but larger specimens, up to six metres (20 feet), have been mentioned in tales of early travellers.

The dugong seems to have had a wider distribution in the past, but in recent times it has been restricted to the warm coastal waters of the Indo-Pacific region. Throughout most of its range it is now much depleted in numbers but exists from the Solomon Islands in the east to the head of the Red Sea in the west and from the Philippine Islands and the Persian Gulf in the north to Brisbane, Perth, and Mozambique in the south. It is entirely marine and rarely even enters estuaries. Manatees, however, seem more adaptable and inhabit the coastal, estuarine, and riverine waters of the tropical and subtropical Atlantic region: *T. senegalensis* from West Africa; *T. manatus*, with two subspecies, from the Caribbean; and *T. inunguis*, landlocked in the Amazon Basin.

Natural history. Sirenians tend to live in groups, and dugongs particularly may be highly gregarious at times. An enormous herd, three and a half by one and a half miles in extent, was recorded off Brisbane, Australia, early in the last century. Both genera are totally aquatic, living almost continuously submerged, and this, coupled with their wary and surreptitious habits, makes them exceptionally difficult to study. They are unable to come out of the water, but manatees placed on land nearby may just manage to return to their natural environment by vigorous wriggling of the body. Breathing is by frequent brief visits to the surface, where the tip of the snout is protruded and the nostrils are silently opened by muscular valves. A manatee has been known to last as much as 16½ minutes between taking breaths. Manatees may also be seen lying at the surface, usually with only their backs visible above the water. Even when basking thus, they are easily frightened by sound or disturbance and will silently submerge. On occasion both dugongs and manatees will thrust their heads entirely out of the water, and they sometimes may even disport themselves like cumbersome seals.

Extremely little is known of the breeding biology of sirenians and particularly of the time scale involved. Manatees are known to be long-lived, and dugongs are believed to have similar lifespans. But there is no information about age of maturity, rate of growth, or number of young produced in a lifetime. Mating in manatees has been observed with the animals in very shallow water lying on their sides. Gestation is known to occupy more than 152 days, and the normal birth is a single calf, which receives much maternal care. Very occasionally twin fetuses are found in dugongs. Suckling from the single pair of pectoral mammary glands normally occurs in a horizontal position, but the occasional suckling of a young one held vertically by a flipper may have given rise to the mermaid myth. There is no marked difference in size between the sexes.

Sirenians are totally herbivorous. Whereas Steller's sea cow fed on marine algae, the dugong and manatees feed entirely on green higher plants, a circumstance that limits their distribution to relatively shallow waters. The habit of the dugong is to feed where there is good growth of "dugong grass" (*Zostera*, etc.), the leaves and underwater stems of which are pulled off by the animal's powerful lips. Manatees seem prepared to feed on marine or freshwater plants whether growing on the bottom, floating at the surface, or even growing on banks of rivers up to a foot above water level. It is this catholicity that has led to the experimental use of manatees in Guyana to keep clear irrigation and transport channels in cane fields that otherwise must be cleared by hand. Virtually the sole enemy of both the dugong and manatees is man. There does not seem to be any commercial exploitation of these

Develop-
ment of
pygmy
forms

Aquatic
behaviour

Extinction
by hunting

animals at present, but both are hunted sporadically and locally for food with the use of nets or harpoons. The dugong still has a considerable residual stock off northern Australia and with protective legislation enforced, could have a valuable future. Manatees, because they come into more confined waters, are still in danger of extermination.

Form and function. Sirenians have torpedo-shaped bodies and, like whales, have lost all external trace of hindlimbs. Their tails likewise are flattened horizontally and provide the main propulsive force in swimming. The forelimbs are small and may help in turning and manoeuvring. Though usually somewhat sluggish, sirenians are capable of considerable speeds for short distances.

The main distinguishing characteristics between the dugong and the manatees are that the former have a downturned snout and a forked tail, while the latter have a straight snout and rounded tail. Both have immensely thick, tough skin, which is nearly hairless. In the dugong

Drawing by J. Helmer based on (manatee) photograph courtesy of Field Museum of Natural History, Chicago; (dugong) from G. M. Allen, *Extinct and Vanishing Mammals of the Western Hemisphere*, Cooper Square Publishers

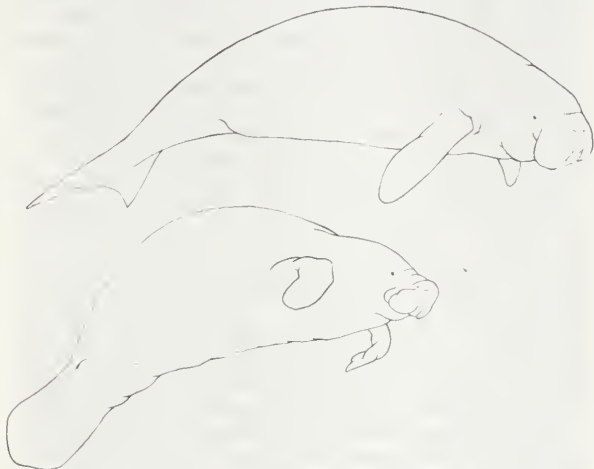


Figure 46: Body plans of two major sirenian groups. (Top) dugong (*Dugong dugon*); (bottom) manatee (*Trichechus manatus*).

there are individual hairs about six millimetres (0.24 inch) long, spaced five to seven centimetres (two or three inches) apart. There is no substantial layer of blubber beneath the skin, but the body contains much fat. The eyes are small and circular, without lids, and the minute external ear openings are to be found only by careful examination. All of the bones are of exceptional density. Sirenians, particularly the dugong, have a much enlarged and intensely muscular upper lip, the corners of which serve to pluck the vegetation on which the animal feeds. In the dugong there are remains of two incisors, the second of which becomes large in the adult, remaining unerupted in the female and erupting to form a tusk in the male. There are basically six cusped molars in each jaw, which fall out progressively from the front so that only two remain in the adult. These teeth, together with horny pads at the front of each jaw, crush the food. In manatees each jaw has a series of 20 to 30 crushing molars, which move progressively forward during the life of the animal. Manatees are also exceptional in having only six vertebrae in the neck, instead of the seven standard in mammals. Almost nothing is known about the physiological processes of these animals.

Classification and paleontology. Sirenians comprise the order Sirenia of the superorder Subungulata, which also includes elephants and hyaxes. Two points of resemblance with elephants are the mode of tooth succession and the pectoral position of the mammary glands. The extensive fossil record indicates that the Sirenia date from the Eocene, about 50,000,000 years ago. More than 20 genera of fossil sirenians have been described. They seem to have been widespread near the coasts of the warmer seas of the world, and some think that the manatee is of more recent origin and has replaced the dugong in the Atlantic region.

(G.C.L.B./C.K.B.)

Perissodactyla (horses, asses, zebras, tapirs, rhinoceroses)

The order Perissodactyla is a group of herbivorous mammals characterized by the possession of either one or three hoofed toes on each hindfoot. The name (Greek *perissos*, "odd," and *daktylos*, "finger") was introduced to separate the odd-toed ungulates from the even-toed ones (Artiodactyla), all of which had previously been classified as members of a single group.

GENERAL FEATURES

The Perissodactyla comprise three families of living mammals: six species of horses (Equidae), four species of tapirs (Tapiridae), and five species of rhinoceroses (Rhinocerotidae). These families are remnants of a group that flourished during the Tertiary Period (from 65,000,000 to 2,500,000 years ago), a time when it was much richer in species and in variety of form than at present and played a dominant role in the fauna of the world.

The horses, asses, and zebras are long-legged, running forms with one functional digit in each foot and with high-crowned, molariform (*i.e.*, modified for grinding) cheek teeth. The tapir is a rather rounded, piglike, semiamphibious forest and woodland animal with a small proboscis (trunklike snout) and a coat of short, bristly hairs. Tapirs have primitive features, such as four hoofed toes in the forefoot and three in the hind, and they have rather simple molar teeth. Rhinoceroses are massive, graviportal (ponderous) creatures with a thick and nearly hairless hide and three digits on each foot. They bear hornlike structures on the head.

The Perissodactyla are of particular scientific interest because their fossil history is so well-known. The evolution of horses from the tiny "dawn horse" (*Hyracotherium*, formerly *Eohippus*) to the present form is a classic sequence, knowledge of which has played an important role in evolutionary thought. The order also provides a notable example of parallel evolution. Following completely different evolutionary paths, both perissodactyls and artiodactyls (*e.g.*, cattle, antelope, swine) independently evolved features such as high-crowned grinding teeth and elongated limbs with a reduced number of digits, in adaptation to a similar running (cursorial), herbivorous mode of life.

Living perissodactyls are of medium or large size. Asses and tapirs, the smallest representatives of the order, attain a length of approximately two to 2.5 metres (6.6 to 8.2 feet), stand one metre or more at the shoulder, and weigh up to 250 or 300 kilograms (550 to 660 pounds). The largest forms are the Indian and square-lipped rhinoceroses (*Rhinoceros unicornis* and *Ceratotherium simum*, respectively), which are four to five metres (13 to 16.4 feet) long and measure up to two metres at the shoulder. The maximum weight has not been well established, but a figure of 2,070 kilograms (4,550 pounds) has been recorded. *Baluchitherium*, relative of the rhinoceroses known from the Oligocene (about 30,000,000 years ago), was the largest known land mammal, standing about 5.5 metres (18 feet) at the shoulder.

Ecologically, the Perissodactyla were the dominant large herbivore group during the Tertiary, a position now held by the Artiodactyla. All feed either by grazing (*i.e.*, cropping grasses) or by browsing (taking shoots and leaves from trees and bushes). The Equidae in particular were abundant and important members of the Old World fauna until their numbers were reduced by modern man. Zebras are still numerous and ecologically important in a few parts of Africa. The importance of the domestic horse and the ass in the history of mankind is very great indeed. Both have served extensively as pack, draft, and riding animals. The horse is sometimes eaten by man, and its flesh is widely used as pet food; through centuries of domestication, it has been developed into a number of different breeds (for more information on domesticated horses, see HORSES AND HORSEMANSHIP).

The living wild Equidae are confined to the Old World. Zebras and the true wild ass (*Equus asinus*) are African, with the zebras confined to the southern and eastern parts,

Indian and square-lipped rhinoceroses

Distinguishing features

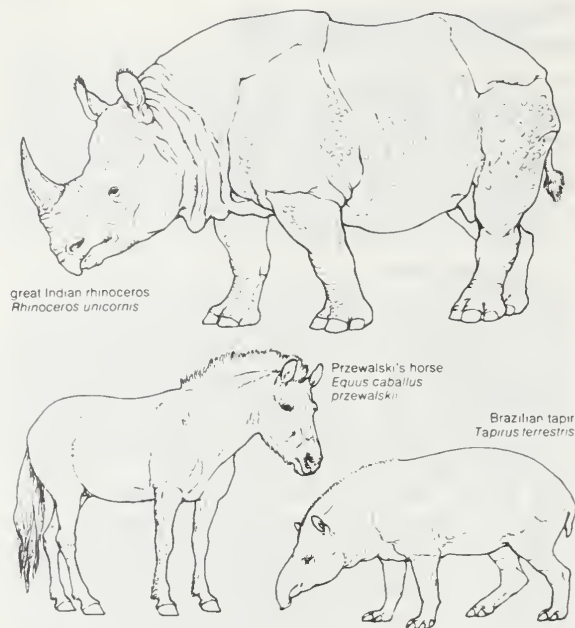


Figure 47: Major body plans among perissodactyls.
Drawing by R. Keane

while the ass originally ranged over northern and north-eastern Africa.

The wild horse (*Equus caballus*), ancestor of the domestic horse, occupied the low country north of the great mountain ranges from Europe across central Asia; it may now be extinct as a wild animal. The half-asses, races of *E. hemionus*, were found in the arid zone of Asia from Persia to the Gobi Desert, as well as in Arabia, Syria, and northwestern India.

The living rhinoceroses are also Old World forms, with two species in Africa and three in Asia. There are three species of tapirs in the New World tropics, one in Middle America and two in South America. The fourth species of tapir is Asiatic.

NATURAL HISTORY

Distribution, ecology, and conservation. The Equidae are highly specialized for a cursorial, herbivorous mode of life. They are absent from forests and other densely vegetated regions, but apart from this limitation, the range of the group is relatively unrestricted by the type of vegetation, climate, and topography. The species replace one another geographically for the most part, and each occupies a somewhat different habitat.

Grass plays a major role in the diet; the zebras, for example, are known to feed on tall, coarse grasses avoided by most antelopes. Some species also take shrubs, herbs, and even bulbs. Water requirements vary in different species. In South Africa, Burchell's zebra has been found to drink about once every 36 hours. By contrast the mountain zebra (*Equus zebra*), Przewalski's horse (*Equus caballus przewalskii*) and the half-ass, all living in semidesert areas, are reported to survive if they can drink once in three or four days. The ass, too, can manage with less water than the horse. The mountain zebra and Przewalski's horse dig for water in dry riverbeds and depressions.

Zebras. The mountain zebra occupies parts of the arid rocky escarpment separating the interior plateau of the southern African subcontinent from the coast lowlands. The race *E. z. hartmannae*, still relatively numerous over much of its original range in South West Africa and southern Angola, enjoys legal protection and is represented in game reserves. Conflict between these zebras and farm livestock for the meagre pasture of the semidesert regions, however, has led to a reduction in zebra numbers. The race *E. z. zebra* was originally common in the mountain ranges of the Cape Province, but now survives only as a remnant of perhaps 100 animals. About one-half of these have sanctuary in a national park.

Burchell's zebra (*E. quagga*) formerly inhabited a great

area of grassland and savanna from the Cape to the southern Sudan. The southernmost race (*E. q. quagga*), which was only partly striped, became extinct in the 19th century. The populations of the other races have been much reduced in many places and the range of the species has shrunk considerably. There are large populations in reserves, however, and the species is not in any immediate danger of extermination. Grevy's zebra (*E. grevyi*), which shares a narrow zone in northern Kenya with Burchell's zebra, is confined to sparsely wooded, semidesert plains and low hills in northern Kenya, southern and eastern Ethiopia, and western Somaliland. Its status appears to be generally satisfactory.

Asses. The true ass (*Equus asinus*), ancestor of the domestic donkey, is the equid of arid North Africa whose range extends south to approximately 6° N latitude. Its natural distribution probably included all habitable parts of North Africa. At present, asses are known from semidesert country extending from the east bank of the Nile (in the Sudan) to the Red Sea, and in parts of Eritrea, Ethiopia, and Somalia. There are also isolated pockets in the Tibesti Mountains in the Sahara, and in the countries of central and western Africa. There is a great deal of uncertainty about the identity of all asses now described as "wild." Some may be merely feral (escaped or released) donkeys, and interbreeding with feral donkeys is likely to have occurred in many, if not all, existing populations.

The wild horse. The wild horse was widely distributed in Eurasia north of the mountain chains. The Romans encountered it in Spain. Two races have survived to modern times. A gray race, known as the tarpan, was the horse of southern Russia. It became extinct in the Ukraine during the mid-19th century. The endangered Przewalski's horse (*E. c. przewalskii*), a small, reddish-brown race (considered a species by some authors), was last seen in the wild in 1968 in the remote, semidesert steppe country on the boundary between Mongolia and China. Wild horses enjoy legal protection in Mongolia and China, but nomadic pastoralists have been encroaching on previously uninhabited country and competing with the horses for pasture and the scarce water supplies.

The half-asses, races of *Equus hemionus*, occupied the dry belt from Mongolia through Central Asia to Syria, with a northern limit at about 50° N latitude. The chigeta or kulan (*E. h. hemionus*), which was formerly widespread over an immense region of the Gobi, now occurs only in semidesert steppe country in central Mongolia. Hunting and competition for water by pastoral tribesmen are responsible for its decline. The kulan is slightly smaller than the kiang (*E. h. kiang*), which is found on the cold, arid steppes of Nepal, Sikkim, and western Tibet at altitudes of 4,270 metres (14,000 feet) and more. The kiang is now said to be rare but not endangered. The Persian onager (*E. h. onager*) lives in a lower semidesert or desert environment, with a range that formerly included northeastern Iran, northwestern Afghanistan, and Russian Turkestan. It is now extremely rare and unlikely to survive outside northeastern Iran and the Badkhyz Reserve in Turkmeniya. A small nucleus has sanctuary in the semidesert salt plains of the recently established Kavir Protected Region in Iran. The Indian wild ass is a closely related, probably identical, form sometimes distinguished as the race *E. h. khur*. A fairly small population occupies salt flats in the Rann of Kutch, a remnant of the thousands found there at the end of World War II. The Syrian onager (*E. hemionus hemippus*) is the smallest member of the group and stands about one metre (three feet) at the shoulder. It was once found in the desert region of Palestine, Syria, and Iraq, and was domesticated by the ancient Sumerians before the introduction of the domestic horse into Mesopotamia. This race may survive in the Djezireh Desert, Syria, or north of the Syrian-Turkish border; if so, the number must be extremely small.

Rhinoceroses. The two African species of rhinoceros are the black or prehensile-lipped rhinoceros (*Diceros bicornis*) and the white or square-lipped rhinoceros (*Ceratotherium simum*). The terms black and white are misleading, since both species are grayish to brownish, but the names are well established in common usage.

The two surviving races of wild horses

Dietary needs as range factors

"White" and "black" rhinos of Africa

The black rhinoceros was originally widespread from the Cape to southwestern Angola and throughout eastern Africa as far as Somalia, parts of Ethiopia and The Sudan. Its range also extended westward through the northern savanna zone to Lake Chad, the northern Cameroons, Volta, and the northern Ivory Coast. The animal was extremely numerous in some parts. It now occupies a much smaller area, within which it is found in scattered pockets, many of them in parks and reserves. South Africa has an important population of some 300 of them in its Zululand reserves. The species still occurs in northern South West Africa and southwestern Angola, Rhodesia, Mozambique, Malawi, Zambia, and Zaire. Northern Tanzania and Kenya have more black rhino than any other country but the future of the animals outside the parks and reserves is far from secure. The species is also found in parts of Somalia, Ethiopia, northern Uganda, and The Sudan, and remnants remain in the Central African Republic and northern Cameroon. The decline in numbers is largely the result of expanding human settlement and of poaching to obtain the horns, which fetch high prices.

The black rhinoceros occupies a variety of habitats, frequenting open plains, sparse thorn scrub, savannas, thickets, and dry forests, as well as mountain forests and moorlands at high altitudes. It is a selective browser, and grass plays a minor role in its diet. Succulent plants, such as euphorbias, assume great importance in dry habitats, and where these plants are abundant the animals appear to be able to survive without free water. Where water is available, drinking is regular and frequent; the animals also may dig for water in dry riverbeds.

The much larger white rhinoceros is a grazing species with a broad, square muzzle. It prefers short grasses seven to 10 centimetres (about three to four inches) high. The animal makes much use of shade trees for resting and is dependent on surface water. The range of the white rhinoceros is markedly discontinuous. South of the Zambesi River it was once extremely common over a fairly large area of bushveld. It has since become confined to the game reserves in Zululand, where the population has risen to some 1,700; some of the animals have been redistributed to several other parks and reserves in southern Africa.

A northern race formerly inhabited the southern Sudan and adjacent areas of Uganda and Zaire, extending westward into the Central African Republic. It has also been much reduced, but considerable populations still survive in the Parc National de la Garamba (Zaire) and in the Bahr el Ghazal region of The Sudan, and a small remnant has found sanctuary in the Ruwenzori National Park, Uganda.

The smallest of the three Asian rhinoceroses (also the smallest living member of the family) is the Sumatran, or Asiatic, two-horned rhinoceros, *Didermoceros* (or *Dicerorhinus*) *sumatrensis*, standing one to one and a half metres (three to five feet) at the shoulder. It was originally found from eastern Pakistan and Assam throughout Burma, much of Thailand, Indochina (Cambodia, Laos, and Viet Nam), Malaya, Sumatra, and Borneo. Small isolated populations still occur in a few widely separated localities in Burma, Thailand, Malaya, Sumatra, and Sabah, and possibly in other nearby territories. The total population is thought to number between 100 and 170. Some of the survivors in Sumatra are protected in reserves.

The Javan, or lesser one-horned, rhinoceros (*Rhinoceros sondaicus*) occupied the islands of Java, Borneo, and Sumatra, the Malay Peninsula, and a region extending northwards through Burma into Assam and eastern Bengal. It is now restricted to the Ujung-Kulon Reserve in western Java where there are at least 25 and perhaps as many as 50 to 60 animals.

Both the Sumatran and Javan rhinoceroses inhabit forests as well as marshy areas and regions of thick bush and bamboo, climbing actively in mountainous country. They are mainly browsers.

The great Indian, or one-horned, rhinoceros (*Rhinoceros unicornis*) is more or less equivalent in size to the square-lipped rhinoceros and is distinguishable from the smaller Javan rhinoceros by the presence of a large horn, tubercles on its skin, and a different arrangement of skinfolds.

It previously occupied an extensive range across northern India and Nepal from Assam in the east to the Indus Valley in the west. It is found in a range of habitats—open grassland, savanna, forests, and hilly country—and appears to be mainly a grazer, raiding grain fields in some areas. Hunting and the pressure of expanding human populations have greatly reduced both the range and numbers of this animal. It is now found almost entirely in eight reserves or sanctuaries in India, notably the Kaziranga Sanctuary in Assam (estimated population 300) and in the Rapti Valley region of the Nepal Terai. The total population is estimated at about 600 animals, and the prospects for survival appear to be reasonably good.

Tapirs. The Malayan tapir (*Tapirus indicus*), largest member of the family Tapiridae, is found in Sumatra and the Malay Peninsula, as far north as the Burmo-Siamese border in latitude 18° N. It is found from sea level to high altitudes and occupies forests and thickets but may feed in more open areas. It is still abundant and widespread.

The three New World species occupy distinct, nonoverlapping but contiguous ranges. The mountain tapir (*Tapirus pinchaque*), the smallest and most primitive, inhabits the temperate-zone forests and bordering grasslands of the Andes in Colombia and Ecuador and in northern Peru, up to altitudes of nearly 4,600 metres (about 15,000 feet). Agricultural and pastoral expansion have resulted in some decline in the status of this species, but it is still fairly common. The Central American, or Baird's, tapir (*T. bairdii*) is the largest of the American species. It is essentially middle American, with a range extending from Mexico into coastal Ecuador, and it occupies undisturbed climax rain forest. It is shy and adjusts poorly to the disturbance caused by settlement. This disturbance, together with the destruction of habitat accompanying human occupation, has greatly reduced its range and numbers. The species is said to be much in need of active conservation. The most widespread species is the Brazilian tapir (*T. terrestris*), which is found throughout the Brazilian subregion east of the Andes and in a small area west of the Andes in northwestern Venezuela and northern Columbia. Like the other species, it is largely a forest form requiring the proximity of water. The three New World tapirs are mainly browsers and are remarkably similar in habits.

Behaviour. Expression and communication. The Equidae communicate by means of calls and changes in facial expression. Six different sounds are made by Burchell's zebra. A whinny, consisting of a series of two- or three-syllabic "ha" sounds, serves to maintain contact between members of a group. The repertoire includes an alarm call ("i-ha"), an alarm snort, a drawn-out snort of satisfaction, and a squeal of pain and fear. Other species utter similar sounds, the whinny of the horse and the bray of the ass being well-known examples. Characteristic facial expressions have been described for greeting ceremonies (mouth open, ears up), threat (mouth open, ears back), and submission (mouth open, nibbling movements, ears down). In all species studied except the horse, females assume a particular expression ("mating face") when permitting the male to mount.

In the rhinoceroses and tapirs snorting, squealing, belching, and, in some forms, whistling sounds play a major role in communication. Visual signals are not well developed in these nonsocial animals, but a few facial expressions are used.

Feeding. The Perissodactyla are mainly grazers or browsers. The quality and quantity of grasses available to grazing species may vary considerably with the season and the area. The animals may accordingly move great distances to reach attractive sources of food. Migrations of Burchell's zebras to succulent pastures during the rainy season are a feature of the Serengeti Plains and the Etosha National Park in Africa. The distribution of asses, half-asses, and horses inhabiting arid areas largely follows that of rainfall and pasture.

The food of the browsers is fairly readily available throughout the year; thus, species in this category are relatively sedentary. The browsing rhinoceroses may break down trees and shrubs, and use their forelimbs to help get at otherwise inaccessible leaves and twigs. Food is plucked

Distribution
of New
World
tapirs

The
smallest
living
rhinoceros

with the lips. In the tapirs, the upper lip is fused with the short proboscis. The rhinoceroses (excluding the white) have a pointed upper lip with a fingerlike process that is used to pluck leaves and twigs. The white rhinoceros, with its broad square muzzle, is the most specialized grazing rhinoceros, feeding on grass.

Social organization and territory. In both the mountain and Burchell's zebras the family group is the basic social unit. It generally consists of a single adult male and two or three adult females with their foals. The groups are stable, apparently because of strong mutual ties among the females rather than because of herding by the male. The stallion is dominant, and there is a hierarchy among the mares, the highest ranking (alpha) animal usually leading the group. Other males are either solitary or live in bachelor groups of two or three, sometimes up to 10. Juveniles leave the family group when they attain sexual maturity at one and one-half to two years. When large aggregations occur on favoured grazing grounds, the groups retain their identity. Among Grevy's zebras and wild asses, territorial males and groups of mares and foals and of stallions are found. There is no evidence for territorial behaviour among any of the zebras except Grevy's. Individual groups occupy home ranges that overlap to some degree with those of other groups.

The social organization of other equids is not as well documented. Observers studying the wild horse and half-asses have noted that females and juveniles form a group dominated by a single stallion, which keeps them together by active herding; the unattached males are solitary or live in small herds.

The pattern of social organization among the rhinoceroses is quite different from that of the Equidae. Dominant adult males of the white rhinoceroses occupy territories that, in the Natal reserves, average about 200 hectares (500 acres). Within its area a male may tolerate subadult or aged bulls, which have subordinate status. Adult females accompanied by their calves inhabit home ranges encompassing the territories of six or seven dominant bulls. Juveniles consort with other juveniles or with calfless females, but groups of more than two usually do not stay together long.

The black rhinoceros is basically solitary. Adults of both sexes usually occupy home ranges of 1,000 hectares (2,500 acres) or more, the size depending on the characteristics of the environment; occasionally the ranges may be as small as 200 hectares, however. There is a good deal of overlap in the utilization of these home ranges.

The Asiatic rhinoceroses also are essentially solitary, but detailed information on the nature of the areas they inhabit is not available. Individual great Indian rhinoceroses are said to occupy tracts as small as eight to 20 hectares (20 to 50 acres).

Little is known about the social organization or territorial behaviour of the tapirs. All species are reported to be found alone or in pairs.

Male zebras and horses follow mares in estrus. The stallion, after smelling the spot where a mare has urinated or defecated, exhibits "flehmen" (a characteristic display in which the head is lifted and the upper lip raised) and then urinates or defecates on the same spot. In similar fashion, members of stallion groups often urinate or defecate consecutively; communal dung heaps formed by five to eight animals often arise in this way. The significance of such behaviour is not clear.

Among the Rhinocerotidae excretory products play an important role in marking territories and home ranges. Dominant male white rhinoceroses defecate almost entirely on heaps within their territories. They then scatter the material by kicking vigorously, presumably leaving an individual scent mark in this manner. In addition, they urinate in a ritualized fashion, spraying the urine in powerful jets in a manner peculiar to them and shown by no other sex or age group. Other members of the population also use dung heaps (either in territories or in communal areas, such as along paths) but not exclusively, and they do not scatter dung.

In the black rhinoceros, dung-scattering behaviour does not appear to be exclusive to dominant males. The function of the communal heaps may be mainly to establish

the presence of the inhabitant in his home range, and to maintain contact between known animals.

Dung heaps and urine spraying are also observed among other species of rhinoceroses and among tapirs; their significance is presumably of a similar nature.

Fighting. The pattern of fighting is related to the amount of lethal equipment the various groups possess. The Equidae, unarmoured, do not employ stylized fighting techniques to reduce the danger of serious injury—as among certain other species. Fighting is largely confined to adult males competing for estrous mares. Various techniques occur in the zebras, which may serve as an example of the family. Circling, neck fighting, biting (either in a standing or sitting position), rearing combined with biting and kicking, and kicking on the run all are used, either alone or in combination. No set pattern is followed.

Fights among rhinoceroses consist of charges and striking with the horns, usually accompanied by vocal threats. Goring is not common, the stylized pattern having probably been evolved to minimize the danger of serious injury from the formidable horns.

Amicable behaviour. Mutual grooming is well known among horses. Two animals stand facing in opposite directions and groom each other by nibbling at the root of the tail and the base of the neck. Burchell's zebra behaves similarly and so, presumably, do other members of the family.

Zebras greet each other simply by nose-to-nose contact, except that adult stallions go through a ceremony involving nose-to-genital contact. Nose-to-nose greeting is also characteristic of tapirs and rhinoceroses. The latter also rub their bodies together.

Courtship and mating. Courtship is relatively simple among the social equids. The true ass is apparently exceptional. The partners are strangers when the first approaches are made and the female requires violent subjugation by the male, which bites, kicks, and chases her before she will stand for him. This may be the result of separation of the sexes outside the mating season. The wild horse and Burchell's zebra are not at all violent. The stallion often grooms the mare before attempting to mount. The estrous mare (especially, or exclusively, young mares in the case of Burchell's zebra) adopts a typical posture with legs slightly apart, tail lifted, and, except in the horse, a characteristic facial expression (the "mating face" already mentioned).

The more or less solitary rhinoceroses and tapirs go through a more elaborate courtship, presumably because the partners are strangers. After a chase, the male and female may engage in low-intensity fighting, ending with the male laying his head on the female's rump and then mounting and copulating for an extended period. Several males may mate with an estrous female.

Relations between parent and offspring. The perissodactyls bear well-developed (precocial) young, usually a single offspring. After the mother has assisted in removing the placenta and has licked her offspring in the usual mammalian fashion, the young animal soon attempts to stand. A Burchell's zebra foal has been observed to stand quite firmly 14 minutes after birth, and a black rhinoceros calf 25 minutes after birth.

Newly born equids follow any nearby object during the first few days of life. At this time, zebra mares drive away all other zebras from their foals. The behaviour ensures that the foal will form a bond with its mother during the initial period of imprinting. Foals follow their mothers closely and are groomed frequently.

Although precocial, black rhinoceros calves appear to have a lying-out period; that is, an initial period when they rest quietly in thick cover except when being suckled. Thereafter they follow their mothers closely. A young white rhinoceros tends to walk ahead of the mother and may be guided by her horn.

Most young perissodactyls remain with their mothers until the next offspring is born. A young rhinoceros may, therefore, accompany its mother until it is two and one-half years old, or older. Although grazing starts early, suckling proceeds for a considerable time, perhaps for its psychological rather than its physiological value.

Herd
behaviour
among
equids

Mutual
grooming
and
greetings

Scent
marking
and
territories

Running
games
and mock
fighting

Play. As in other mammals, play is a prominent form of behaviour among young perissodactyls. Zebras up to the age of one year frequently engage in running games. Foals gallop wildly about on their own, jumping and kicking up their heels, sometimes chasing other animals, such as gazelles, mongooses, or birds. In groups they play catching games, running after one another in close succession. Mock fighting sometimes takes place. Groups of adults have also been seen to chase foal groups in play, and indeed stallion groups carry out playful gallops. Stallions also engage in play greeting and in mock fights.

Playful romping and mock fighting with the horns are common among rhinoceros calves. Young tapirs play running games.

Rolling and wallowing. Behaviour for the care of the body is widespread among the perissodactyls. Equids frequently roll in dry, loose soil forming rolling hollows—a common feature of zebra country.

Wallowing, which may help regulate body temperature, probably is mainly a form of self-grooming; it is practiced by all species of rhinoceroses. They often spend hours lying in pools during the middle of the day in hot weather. Mud of suitable consistency induces wallowing, which may be followed by sand bathing. Prolonged rubbing on tree trunks or suitable stumps follows a wallow; old rubbing stumps and stones may take on a shine from repeated use.

Tapirs may have the most pronounced tendency to bathe and wallow, but few details of their behaviour are known. They are also said to enter water when disturbed.

Reproduction. Female equids of all the species for which information is available attain puberty at about one year, but are not normally successfully mated before the age of two to two and one-half years, and possibly as late as three to four years in the case of Grevy's zebra. Zebras probably breed until about 20 years of age. The domestic species are seasonally polyestrous (repeatedly fertile), coming into breeding condition in spring and, unless mated, undergoing repeated estrous cycles at intervals of approximately three weeks until the end of the summer. The wild species studied also tend to mate seasonally; most young are born in spring and summer.

Repro-
ductive
potential
of equids

The gestation period of equids is between 11 and 12 months. In most species a postpartum estrus occurs, usually within two weeks of the birth of the young; thus, the maximal potential reproductive rate is one young per year. This potential is not always attained. Only about 50 percent of domestic mares that are mated produce foals, and nearly half of a study group of Burchell's zebra mares bore only one foal in three years.

The gestation period of three species of rhinoceroses is about 15 to 17 months. For the Sumatran rhinoceros the period is said to be only seven months. No information is available for the Javan rhinoceros. Female white and black rhinoceroses attain sexual maturity at the age of four to five years and are capable of calving at intervals of approximately 2½ years. Rhinoceroses probably breed until between 30 and 40 years old. The white tends to have a mating peak in spring, corresponding with the flush of green grass, and a calving peak in autumn.

The Malayan and Brazilian tapirs have gestation periods of 13 months' duration. The Brazilian tapir is reported to mate before the onset of the rainy season.

FORM AND FUNCTION

Integument. The skin of rhinoceroses is extraordinarily thick. The great Indian and Javan rhinoceroses are covered with large, practically immovable plates, separated by joints of thinner skin to permit movement. The two species differ in the arrangement of the folds. The hair of all rhinoceroses is sparse or absent except that of young Sumatran rhinoceroses, which have a dense coat of crisp, black hair. The skin of the tapirs is also thick with a sparse covering of short hairs arranged in irregular groups. The equids have a normal hide with a well-developed hairy coat.

Structure
of
rhinoceros
horns

The "horns" of the rhinoceroses are noteworthy structures of epidermal origin. The horn is composed of a mass of fused epidermal cells that are impregnated with a

tough, fibrous protein (keratin) and that rest on a roughened bony cushion on the fused nasal bones. The male Javan rhinoceros has a short horn about 25 centimetres (10 inches) long; the horn of the female is rudimentary. The great Indian rhinoceros has a single horn up to 60 centimetres (25 inches) long. The other species of rhinoceroses have a second horn that stands on a protuberance of the bones (frontals) between the eyes. Hornlike structures were also present in the titanotheres and extinct rhinoceroses.

In all living perissodactyls the terminal digital bones are flattened and triangular, with evenly rounded free edges, and are encased by keratinous hooves derived from the integument. The single hoof of the equids—the only mammals to walk on the tips of single digits—is the most highly developed structure of this kind among mammals. The keratinous wall is analogous to the nail of mammals that have claws or nails.

Skeleton. Backbone. The vertebral column acts as a firm girder, with high dorsal (neural) spines on the thoracic vertebrae, above the forelimbs and ribs. Spines and ribs serve as compression struts above and below. The column balances largely on the forelegs and is pushed from behind by the hindlegs, which are the main propellants. This skeletal structure permits running and also enables great weights to be borne in such animals as the rhinoceroses. There are never fewer than 22 thoracolumbar (trunk) vertebrae.

The neck, or cervical, vertebrae are opisthocoelous; *i.e.*, with the bodies (centra) of the vertebrae hollowed behind to take the convex heads of the succeeding centra. This feature facilitates rotatory movement of the neck and is most highly developed in the horses.

Limb girdles. The shoulder blade is long and narrow with a small coracoid process (a ridge to which muscles are attached) and a low spine. There is no clavicle (collarbone). The pelvic girdle has a broad, vertically raised ilium to which are attached the large gluteal (thigh) muscles, important for locomotion, and the abdominal muscles, which carry the weight of the belly.

Limbs. There is a clear evolutionary tendency in the Equidae for the limbs to become long and slender, with a reduction in the number of digits in the swift-running forms. These changes are accompanied by an increase in rigidity and specialization for movement fore and aft. The upper (proximal) segments, the humerus of the forelimb and the femur of the hindlimb, have remained short. In contrast the lower (distal) parts, consisting of an anterior radius and posterior ulna in the forelimb and an anterior tibia and posterior fibula in the hindlimb, have become longer and slimmer. The humerus is short and broad. Its articulation with the radius and ulna only permits forward-and-aft movement. The proximal end of the stout femur has a third projection, or trochanter, in addition to the usual mammalian two, which serves as an additional point of attachment for the large locomotory muscles of the hindlimbs.

Evolution
of swift-
running
forms

The anatomical feature of the order now considered most significant is that the axis of symmetry of the limbs passes through the third or middle toe, the most strongly developed and the one on which most of the weight is borne. This is called the mesaxonic condition and is contrasted with the paraxonic condition of the Artiodactyla, in which the axis passes between the third and fourth toes.

Originally the five toes of the limb were held in the semidigitigrade position—*i.e.*, with the weight of the body being borne on the soles of the toes and on the lower ends of the elongate metacarpal and metatarsal bones of the forefeet and hindfeet. The upper (proximal) ends of these bones were raised above the ground, a condition still to be seen in the tapirs.

As the third digit became increasingly dominant, it became longer and thicker. The upper ends of the third metacarpus and metatarsus broadened and forced the other digits to the side. The first (inner) digit was the first to disappear. The earliest known forms already bore only three digits on the hindfoot. The loss of the first toe on the front foot led to the four-toed condition common in the Eocene; the fifth digit persisted, although it was somewhat

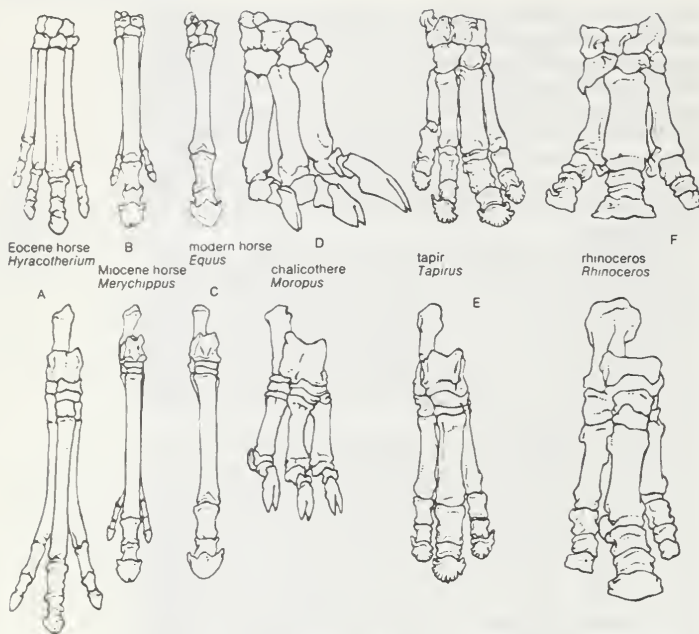


Figure 48: Forelimbs (above) and hindlimbs (below) of some recent and fossil perissodactyls.

Drawing by R. Keane based on (A, B, C, D) A. S. Romer, *Vertebrate Paleontology*, Copyright 1933, 1945, 1966 by the University of Chicago, all rights reserved.

weak. The tapirs, with four toes in front and three behind, have retained this early condition.

The fifth digit of the hindfoot was the next to disappear in the evolutionary sequence, and all known forms beyond the lower Oligocene had only three toes in each foot. The living rhinoceroses illustrate this condition. A massive limb with a broad foot is essential in such heavy animals. The feet have cushions of elastic connective tissue well suited to bear the weight of the body.

In the most highly specialized forms, the second and fourth digits also underwent reduction. These digits are retained in the living Equidae only as functionless, vestigial slivers of bone on either side of the third metacarpal and metatarsal.

With reduction in the number of digits, the third has assumed a progressively more vertical position. Its terminal joint, or phalanx, has become larger and the hoof surrounding it bigger and thicker. At the same time the ulna, the smaller bone of the forelimb, decreased in size until, in the modern Equidae, its upper end fused with the radius and its lower end remained merely as part of the articulating surface of the radius with the wrist bones (carpals). In the hindlimb, the fibula became reduced in similar fashion. It articulates with the ankle bones (tarsals) only in a few extinct forms, such as the titanotheres. In the tapirs and rhinoceroses it is slender. In the equids the proximal end remains as a small splint of bone, while the distal end has fused with the tibia.

The arrangement of the small bones of the carpus ("wrist") and tarsus (ankle) is another characteristic feature of perissodactyl limbs. In most other mammalian orders, the carpals and tarsals provide the limbs with regions of flexibility. Their development in the perissodactyls has been toward increasing rigidity, following the trend in the limbs as a whole. The basic mammalian arrangement is one of three rows of bones: a proximal row of three, adjacent to the lower arm and leg bones; an intermediate row of four (the centralia); and a distal row of five (carpalia of the forefoot, tarsalia of the hind), one to each digit. The centralia, carpalia, and tarsalia are numbered one to four or five, beginning on the side of the thumb and big toe. In all orders the number of wrist and ankle bones has been reduced. Usually only one of the centralia, known as the navicular, remains, while there are the same number of carpalia and tarsalia as there are digits.

An interlocking plan is characteristic of the Perissodactyla. In the forefoot, the third distal carpal (the magnum or capitae) is enlarged and interlocks with the proximal

carpals. The elongated third metacarpal thrusts up against these interlocked bones. In the equids, distal carpal I (the trapezium) is absent, and the arrangement in the hindfoot is similar. In the most advanced, modern forms metatarsal III thrusts against the enlarged and flattened tarsal III (ectotuneiform), and this in turn is in contact with the large, flat navicular (centrale II and III). The navicular abuts on the flattened astragalus (or talus), the intermediate bone of the proximal row. The articulation between the upper surface of the astragalus and the tibia is pulley-like and permits only fore-and-aft movement of the limb.

Teeth. The full complement of mammalian teeth consists of three incisors, one canine, four premolars, and three molars in each half of each jaw. The arrangement

may be expressed by the formula $\frac{3 \cdot 1 \cdot 4 \cdot 3}{3 \cdot 1 \cdot 4 \cdot 3} = 44$ teeth.

The figures represent the number of incisors, canines, premolars, and molars in each half of the upper (above the line) and lower (below) jaws, respectively.

The Condylarthra, a group of mammals that first appeared in the Paleocene (about 65,000,000 years ago) and were ancestral to most of the later and recent hoofed mammals, had a full complement of teeth. In many of the early perissodactyls, only the first lower premolar had been lost. The subsequent evolutionary sequence led to losses and specializations of the incisors and canines. Lengthening of the facial part of the skull resulted in the formation of a gap, the diastema, between the incisors and the premolars. The first upper premolar was reduced or lost in consequence. The canines, when present, were situated in this diastema, as they are in male horses.

Among the living perissodactyls, the tapirs have the least specialized battery of teeth. In this, as in many other features, they have remained primitive. The dental formula of the family Tapiridae is:

$$\frac{3 \cdot 1 \cdot 4 \cdot 3}{(2-3) \cdot 1 \cdot 3 \cdot 3} = 40-42$$

teeth. The first upper premolar is noteworthy in being the only premolar with a milk, or deciduous, predecessor.

The cutting teeth are reduced in the Rhinocerotidae. Incisors and canines are absent in the two African forms. The Asiatic species have one or two upper, but generally no lower, incisors in each half of the jaw. They have no upper canines; the Javan and Sumatran rhinoceroses have one short, sharp lower canine on each side. There are three upper and lower molar teeth in all five species. The white and Sumatran rhinoceroses have three premolars, and the others have three or four premolars. The dental formula for the family is set up,

$$\text{therefore, as follows: } \frac{(0-2) \cdot 0 \cdot (3-4) \cdot 3}{0 \cdot (0-1) \cdot (3-4) \cdot 3} = 24-30 \text{ teeth.}$$

The dental formula of the Equidae is $\frac{3 \cdot 1 \cdot (3-4) \cdot 3}{3 \cdot 1 \cdot 3 \cdot 3} = 40-42$ teeth.

The form of the premolars and molars is of great interest and their evolutionary history has been studied in some detail. Primitive, browsing members of the order had brachydont cheek teeth (*i.e.*, with low crowns and long, narrow root canals), with separate low, rounded cusps—the bunodont condition. Increasing specialization for grazing resulted in fusion of the cusps into ridges (lophs), thus teeth of this kind are called lophodont. Lower molars typically have two transverse lophs, the protoloph and the metaloph. In the upper molars these ridges are fused with a longitudinal ridge (ectoloph), which runs along the outer edge of the tooth. Further development leads to a convoluted arrangement of the lophs, such teeth being termed selenolophodont.

Associated with these changes in the tooth surfaces is a tendency for the crown to become higher. High-crowned teeth are termed hypsodont. The hollows between the lophs of hypsodont teeth are filled with a deposit of secondary cement, which strengthens the teeth and makes them more resistant to wear. A further evolutionary trend is for premolars to become as large as molars. Where the process of molarization is complete, as in horses, all grinding teeth are identical.

The Tapiridae have primitive brachydont premolars and

Fully-toothed early hoofed mammals

Growth of the third digit's terminal joint

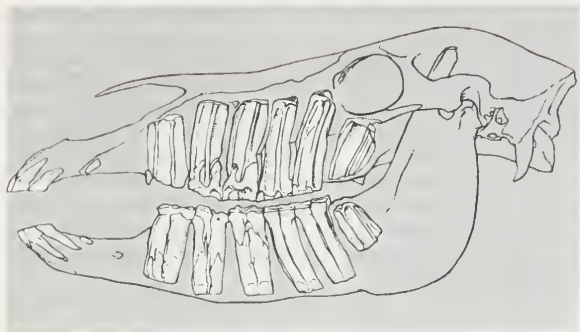


Figure 49: Skull and dentition of horse, *Equus caballus*, cut away to show roots of teeth.

From *Horses. The Story of the Horse Family in the Modern World and Through Sixty Million Years of History*, by G. G. Simpson. Copyright 1951 by Oxford University Press, Inc. Reprinted by permission.

molars. They possess two simple transverse lophs, and are thus termed bilophodont. They lack secondary cement on the crowns.

The Sumatran rhinoceros, the most primitive of the living rhinoceroses, and the Javan rhinoceros have similar brachydont, lophodont cheek teeth. The great Indian rhinoceros, which is less of a specialized browser, has hypselodont (hypsodont and selenodont) premolars, with a layer of cement on the crowns. The black rhinoceros has brachydont and lophodont teeth, with a thin layer of cement. The white rhinoceros is more specialized, for the cheek teeth are hypselodont and have a thick cement layer.

The grinding teeth of the Equidae are highly specialized, high crowned, with a complicated selenodont surface and thick cement deposits.

Digestive system. The stomach of perissodactyls is small, simple, and undivided. In the horse its capacity is only 8.5 percent of the whole digestive system. The comparable figure for the ox is 71 percent. The intestine is very long and the cecum (blind gut) and colon are huge and sacculated (*i.e.*, with many blind pockets). Here food is macerated and fermented and the fibrous portions are dissolved. The liver has no gall bladder.

EVOLUTION AND PALEONTOLOGY

The Perissodactyla appeared early in the Eocene, about 55,000,000 to 40,000,000 years ago. Together with most other ungulate mammals, they were probably derived from the Condylarthra. The condylarths were abundant in Europe and North America, mainly during the Paleocene (65,000,000 to 55,000,000 years ago). Condylarths were unspecialized mammals, rather carnivore-like in appearance. The larger species attained the size of tapirs. The limbs were fairly short and primitive but the third toe was somewhat enlarged and the phalanges ended in hooves. There was a full complement of teeth with some specialization of bunodont molars for herbivorous diet.

Horses. The earliest horses appeared during the lower Eocene in Europe and North America. They are generally known as *Eohippus* ("dawn horse"), but *Hyracotherium* is the correct taxonomic designation. Some species of these little forest-dwelling, browsing animals were no larger than a terrier. They had moderately long, slender limbs with only four toes in the forefoot and three in the hindfoot, all equipped with hooves. The molars were essentially bunodont (with low, rounded cusps) and the premolars simple.

Hyracotherium-like animals persisted in Europe until the end of the Eocene. Another group, the paleotheres or "native" European horses, evolved as a specialized side branch, which died out in the Oligocene. North America was the centre of horse evolution. During the Eocene, *Hyracotherium* was succeeded by forms such as *Orohippus* and *Ephippus*, which are known only from that epoch.

The Oligocene (40,000,000 to 26,000,000 years ago) saw a major change with the appearance of three-toed horses, *Mesohippus*, *Miohippus*, and others. All of the premolars were similar to the molars, low-crowned but lophodont (ridged). *Anchitherium* was an early Miocene form as large as a modern pony, which migrated from North America to

Europe. These primitive three-toed horses or anchitheres survived until Pliocene times, some of their descendants attaining the size of a rhinoceros.

The main course of horse evolution entered a third stage in North America in the Miocene Period (26,000,000 to 7,000,000 years ago). A line of grazing horses developed, almost certainly to exploit the new grasslands that were spreading over the surface of the earth. The degree of lophodonty of the molariform teeth increased, changing the pattern of the crests on the surface and increasing their grinding efficiency. Of greater importance, these teeth became hypsodont (high-crowned) and thus maintained a good grinding surface as grinding of the harsh, siliceous grass caused them to wear down. Another substance, cement, came to supplement the dentine and enamel forming the teeth of earlier types, and provided additional material to resist abrasion. The evolution of these specialized teeth was a tremendous advance.

The limbs of the grazing horses became increasingly rigid and specialized for fore-and-aft movement, better fitting the animals for running in open country. In *Merychippus* the ulna was fused with the radius and the fibula was much reduced. In some advanced forms the central toe was much larger than the two lateral toes and carried most of the weight of the body on a hoof much like that of modern horses.

A number of evolutionary lines developed during the Pliocene, which lasted from 7,000,000 to 2,500,000 years ago. *Pliohippus* of North America is probably the line from which modern horses have come. The genus *Equus* is characteristic of the Pleistocene when it developed in North America and spread to all continents except Australia. By the end of the Pleistocene, horses had become extinct in the New World.

Titanotheres. Another group entirely, the titanotheres (*Brontotheriidae*), evolved independently from *Hyracotherium*-like ancestors and became abundant in North America during the Eocene and Oligocene, but disappeared by the Miocene. They were also found in Asia and Eastern Europe. The end forms, such as *Brontops* and *Brontotherium*, were huge, the largest standing 2½ metres (eight feet) at the shoulder. They had long, low skulls and a small brain. Many species bore a pair of large, hornlike processes on the front of the head.

Chalicotheres. The chalicotheres (*Chalicotheriidae*) were moderately large animals that appeared in Eurasia and North America during the Eocene. Thereafter they evolved mainly in the Old World, disappearing from America in the mid-Miocene, but persisting in Asia and Africa until they died out in the Pleistocene. Early mem-

By courtesy of the Field Museum of Natural History, Chicago



Figure 50: Reconstruction of fossil Perissodactyla. *Brontops*, an Oligocene titanotheres, some individuals of which were about 2.5 metres (eight feet) at the shoulder. Detail of a painting by Charles R. Knight.

Comparative size of the stomach

Trends in horse evolution

bers of the group such as *Paleomoropus*, from the lower Eocene, resembled contemporary equids. The Miocene *Moropus* typifies the peculiar characteristics of later forms. It was horse-like, but the front legs were longer than the rear. Each foot bore three toes which ended in large, fissured phalanges bearing claws that may have been used for digging up roots and bulbs.

Tapirs. The fossil and living tapirs, together with similar extinct forms such as the lophiodonts, constitute a small, fairly uniform group.

Homogalax, *Lophiodon* and other tapir-like animals appeared in the Eocene. Some were close to the ancestry of both rhinoceroses and tapirs, and other evolutionary lines left no descendants. *Protapirus* from the Oligocene of Europe and North America was the forerunner of the modern tapirs. These tapirids persisted until the Pleistocene, when climatic changes led to their extinction.

Rhinoceroses. The fossil history of the rhinoceroses is far more complex. They evolved from the early tapiroids but diverged from that group. Unlike the equids, they have tended to grow large, with short, stout limbs with little reduction in the digits. Although the premolars have tended to molarize, the molar cusp pattern is simple; the teeth seldom become hypsodont, and cement occurs rarely.

The Hyracodontidae, or running rhinoceroses, from the Eocene and Oligocene of North America and Asia were the most primitive. *Hyracodon* had long, slender legs with three toes on each foot and typically rhinocerotid cheek teeth.

The amynodonts are known from the late Eocene and Oligocene of Eurasia and America and lived in Asia until the Miocene. They were a side branch, perhaps derived from primitive hyracodonts. *Metamynodon* and some other forms were about as large as hippopotamuses and may have lived in rivers. The premolars were simple and the incisors reduced, but canines and molars were enlarged.

True rhinoceroses are probably also descended from early hyracodonts. They became numerous in the Oligocene as large animals with molariform premolars. Many side branches evolved. *Trigonias* still had four front toes but *Caenopus*, another Oligocene representative, had the three toes common to all later rhinoceroses. Like its relatives of that period, *Caenopus* was fairly small, about the size of a tapir, and hornless.

One spectacular group contained giant hornless creatures, such as *Baluchitherium* and *Indricotherium* from the Oligocene and Miocene of Asia; the largest of known land mammals, they stood 5.5 metres (about 18 feet) high at the shoulder and had a long neck and long forelegs.

Rhinoceroses died out in North America during the Pliocene, but in Eurasia many different species survived to the Pleistocene. One of the most familiar of these later rhinoceroses is the woolly rhinoceros (*Coelodonta*), which was depicted by Stone Age artists and is known from nearly intact specimens. The surviving rhinoceroses represent remnants of once varied and abundant stock; the genera have little direct relationship to one another.

CLASSIFICATION

Distinguishing taxonomic features. Skeletal features are of greatest importance in classifying the Perissodactyla. They are, of course, the only criteria applicable to fossil forms. Distinguishing characteristics of the skull include the relative length of the facial region, length and form of the nasal bones, presence of a postorbital bar and of hornlike structures. The form and number of teeth, presence of a diastema, number of molariform premolars, the height of the cheek teeth, their cusp structure and the presence or absence of cement on the grinding surfaces all are particularly valuable taxonomic features. The limb structure is also a useful aid to classification, especially the length of upper and lower portions, the degree to which the ulna and the fibula are reduced and fused, respectively, with the radius and the tibia, the number and form of carpal and tarsal bones, and the number and relative size of the digits.

Among living perissodactyls, body size, form of the upper lip, number and length of horns, structure of the skin

and colour pattern are the most important features for classification.

Annotated classification. The classification presented here follows that of U.S. paleontologist George Gaylord Simpson, which is generally accepted but modified in some fairly minor ways by other authors. The term Mesaxonia, introduced by the 19th century paleontologist O.C. Marsh, is essentially synonymous with Perissodactyla. Simpson used it as a convenient designation for the superorder containing the single order Perissodactyla, as he used the parallel name Paraxonia for the superorder with the one order Artiodactyla.

Groups indicated by a dagger (†) are known only as fossils.

ORDER PERISSODACTYLA

Herbivorous ungulates with either 3 digits or 1, at least in the hindfeet. Early forms digitigrade, this feature either retained or completely replaced by the unguigrade condition in later representatives. The body weight borne mainly or entirely on the third digit through which the long axis of the limb passes (mesaxonic condition). The talus (heel bone) with only 1, proximal, keeled surface, articulating with the tibia, no duplication of keel on the distal surface, as in artiodactyls. Nasals broad at the posterior end, alisphenoid canal present. Posterior premolars molariform; cheek teeth bunodont in early forms, typically lophodont or selenolophodont, rectangular; wide diastema separates them from incisors and canines, which may be reduced or absent. Testes inguinal, occasionally scrotal; mammary glands inguinal, with 2 teats; uterus bicornuate. Fifteen living species, but more than 200 fossil forms.

Suborder Hippomorpha

Superfamily Equoidea

Dentition complete, upper molars with 6 tubercles, the 2 external ones united to form an ectoloph, median and internal tubercles generally fused into a single loph. Tendency to molarization of premolars and reduction of lateral digits.

†*Family Palaeotheriidae* (paleotheres or "native horses"). Lower Eocene to lower Oligocene; Europe and Asia. Size variable, from that of fox (shoulder height about 40 centimetres or 16 inches) to that of a rhinoceros; body form tapirlike. Three digits in each foot. Dentition complete or first premolars absent; premolars molariform, separated from canines by diastema. Cheek teeth brachydont, lophodont, with or without cement. Orbits without postorbital bar. Eurasian "native horses." 7 genera.

Family Equidae (horses and relatives). A progressive group with tendency to molarization of premolars and development of cement.

†*Subfamily Hyracotheriinae* (hyracotheres). Eocene; Europe and North America. Size about that of a fox; 4 toes in forefoot, 3 in hindfoot. Dentition complete, premolars simple or last 2 slightly molariform, molars brachydont, bunodont, diastema short. Skull short, without postorbital bar. Three genera. Basic stock from which other members of family evolved.

†*Subfamily Anchitheriinae* (anchitheres). Lower Oligocene to lower Pliocene. Fossils from North America, Europe, Asia. Browsing horses. Early forms collie-sized (shoulder height about 60 centimetres [23.6 inches]), end forms as large as a rhinoceros. Three toes in each foot, central toe somewhat enlarged; metapodials elongated. First premolars reduced, second to fourth molariform. Cheek teeth brachydont, lophodont; ectoloph (longitudinal outer crest) above typically W-shaped cross-lophs retained distinct cusps; hypostyle at back edge of upper molars prominent. Lower cheek teeth with 2 crescent-shaped lochs and double cusp at junction. Diastema long, no postorbital bar. Six genera.

Subfamily Equinae (horses, asses, half-asses, and zebras). Lower Miocene to Pleistocene of North America; Pleistocene of South America; lower Pliocene to Recent of Europe, Asia, and Africa. Most forms as large as ponies or modern horses. Feet with 3 toes in early forms, 1 in later forms (but metapodia 2 and 4 persist as remnants), limbs with progressively elongated lower segments and reduced or fused ulna and fibula. Calcaneum without articulating surface for fibula. First premolars reduced or absent, others molariform; cheek teeth hypsodont, selenolophodont, with cement; canines in long diastema. Skull long, with postorbital bar in upper Miocene and later forms; nasals long and narrow. Some Recent forms with striped coat.

†*Superfamily Brontotherioidea*

†*Family Brontotheriidae* (titanotheres). Lower Eocene to lower Oligocene. Fossils from North America and Asia. Early forms as large as tapirs, end forms huge, standing 2.5 metres (8.2 feet) at the shoulder. Graviportal browsers; 4 toes in forefoot, 3 in hindfoot. Dentition complete, premolars small

and simple in early forms; later forms with incisors and first premolars reduced or absent, remaining premolars only partly molarized. Molars large, bunolophodont, above with W-shaped ectoloph as in early equids, lower with double V resembling equid pattern but lacking reduplication of cusp at union of V's. Skull elongated, low, without postorbital bar, with large rough, hornlike process in later forms. About 39 genera.

†Suborder Ancylopoda

†Family *Chalicotheriidae* (chalicotheres). Upper Eocene to lower Pliocene. Fossils from North America, Europe, Asia. Early forms equivalent in size (and similar) to contemporary horses, *Hyracotherium*; later representatives larger, up to size of a modern horse. Most with forelegs longer than hindlegs and 3 toes in each foot; rudimentary metacarpal V retained in forefoot, most with clawlike digits. Incisors of canines reduced or absent; first premolar absent, remainder small and simple. Molars bunolophodont, cusp pattern similar to Brontotheriidae. About 13 genera. Aberrant perissodactyls which may have used claws to dig for bulbs and roots.

Suborder Ceratomorpha

Superfamily Tapiroidea

Brachydont forms, molars with simple ectoloph and strongly developed transverse protoloph and metaloph. Forefoot generally with 4 toes, hindfoot with 3.

†Family *Isectolophidae*. Lower to upper Eocene. Fossils from North America and Asia. Limb structure similar to that of *Hyracotherium*, teeth already tapiroid. Four genera.

†Family *Helatidae*. Lower Eocene to middle Oligocene. Fossils from North America and Asia. Skull advanced in structure; premolars much simpler than molars, molars of simple lophodont form. Ten to 12 genera.

†Family *Lophiodontidae* (lophiodonts). Lower to upper Eocene; Europe and probably Asia. Skull not tapir-like, lophs of molariform teeth oblique. About 6 genera.

Family *Tapiridae* (tapirs). Lower Eocene to Recent. Fossils from North and South America, Europe, Asia, Africa. Recent species in Central and South America and Southeast Asia. Limited to tropical equatorial forests of South America and Asia. Browsers. Skull greatly modified; premolars completely molarized, first sometimes absent; transverse lophs strongly developed in molariform teeth, ectoloph simple. Limb structure primitive; 4 toes in forefoot; 3 in hindfoot. Six genera but only 1 extant.

Superfamily Rhinocerotoidae

Molars brachydont, upper with 2 transverse lophs fused with a well-developed ectoloph; lower with 2 asymmetrical crescents. Premolars more or less molariform. Rough processes on nasals (for horns) frequently present; forefoot with 3 or 4 toes.

†Family *Hyrachyidae*. Lower to upper Eocene; North America and probably Asia. Small animals; forefoot 4-toed; molariform teeth with separate tubercles. Four genera.

†Family *Hyracodontidae* (hyracodonts or running rhinoceroses). Middle Eocene to upper Oligocene. Fossils from North America and probably Asia. Relatively small; feet 3-toed, legs long and slender. Molariform teeth typically rhinocerotid. About 7 genera.

†Family *Amynodontidae* (amynodonts). Upper Eocene to middle Oligocene. Fossils from North America, Europe, Asia. Hippopotamus-like; forefoot 4-toed, hindfoot with 3 toes. Skull short, heavy, premolars and incisors reduced, premolars simple, canines and molars enlarged. Six genera.

Family *Rhinocerotidae* (rhinoceroses). Middle Eocene to present. Fossils from North America, Europe, Asia, and Africa. Now restricted to tropical Africa and Asia. Large; most forms 3-toed. Skull elongated, raised at posterior end. Dentition incomplete, upper canines always absent, some incisors usually enlarged as cutting teeth; premolars molariform, last upper molar simple. About 34 genera, 4 extant.

Critical appraisal. Systems of classification frequently differ in the status given to the chalicotheres. Simpson takes the view that their ancestry was probably equid and almost certainly hippomorph. He holds that the significance of their claws has been over-emphasized and has tended to distract attention from their true affinities. Accordingly he places them in the superfamily Chalicotheroidea of the suborder Hippomorpha. In the classification presented above they are given their own suborder, Ancylopoda, following the views of Alfred S. Romer, another authority on the group.

The forms united in the family Lophiodontidae by Simpson, followed here, are thought by some recent workers

to warrant separation into three families. Romer distinguishes the families Lophialetidae, Deperetellidae, and Lophiodontidae. The affinities of certain primitive genera such as *Hyrachus* and *Colonoceras* remain controversial; Simpson places them in the family Hyrachyidae, superfamily Rhinocerotoidae. Romer considers them to be early tapiroids and assigns them to the Helatidae (superfamily Tapiroidea). The difference of opinion is slight, for it is generally agreed that the hyrachids are close in the common stem of the tapiroids and rhinocerotoids.

(R.C.Bi.)

Artiodactyla (pigs, goats, sheep, cattle, giraffes, deer, camels)

The mammalian order Artiodactyla, or even-toed ungulates, includes the pigs, peccaries, hippopotamuses, camels, chevrotains, deer, giraffes, pronghorn, antelopes, sheep, goats, and cattle. It is one of the larger mammal orders, containing about 150 species, a total that may be somewhat reduced with continuing revision of their classification. Many artiodactyls are well-known to man, and the order as a whole is of more economic and cultural importance than any other group of mammals. The much larger order of rodents (Rodentia) affects man primarily in a negative way, by competing with him or impeding his economic and cultural progress.

GENERAL FEATURES

Abundance and distribution. Artiodactyls were once the dominant herbivores (plant-eating mammals) of almost every continent. They are an important link in the chain by which the sun's energy, having been used by green plants, is made available to other forms of life. They tend to be medium- or large-sized animals. If they were any smaller they would compete with rabbits and the larger rodents, and if they were larger they would compete with elephants and rhinoceroses, the largest of terrestrial herbivores. The success of artiodactyls has depended on skeletal adaptations for running and on the development of digestive mechanisms capable of dealing with plant foods; none is adapted to flying, burrowing, or swimming. The individual species tend to be fairly narrowly adapted, in comparison with other mammals, but many of them nonetheless have broad distributions.

Native artiodactyls are absent only from the polar regions and from Australasia, but many have been introduced into Australia and New Zealand. In Australia, the position of medium and large herbivores is occupied by kangaroos. Through most of its evolutionary history, the order was absent from South America; only within the last few million years have some groups entered that continent. The occurrence of the majority of living artiodactyls in the Old World is a recent phenomenon; a considerable variety once inhabited North America.

The order Artiodactyla contains nine families of living mammals, of which the Bovidae (antelopes, cattle, sheep, and goats) is by far the largest, containing nearly 100 species. There are five Eurasian and four African species of pigs (family Suidae) and two Central and South American species of piglike peccaries (Tayassuidae). The two hippopotamus species (Hippopotamidae) are African. The more familiar large species were until recently widespread throughout Africa south of the Sahara and in the Nile Valley; the pygmy hippopotamus has a restricted distribution in West Africa. The camel group (Camelidae) was formerly abundant in North America, the now extinct North American stocks having produced the camelids of South America (wild guanaco and vicuña, domestic llama and alpaca) and the Old World dromedary and Bactrian camel.

The remaining artiodactyls (*i.e.*, the suborder Ruminantia) are all ruminants (cud chews), the most primitive of which are the chevrotains (Tragulidae), with three species in Asia and one, the water chevrotain, in West Africa; the chevrotains are clearly remnants of a group that was once more numerous and widespread. Deer (Cervidae) are basically Eurasian and have not spread into sub-Saharan Africa, although they have reached the Americas. There are about 30 species, the greatest number being concen-

Families of living artiodactyls



Figure 51: Representative non-bovid artiodactyls.
Drawing by R. Keane

trated in South America and tropical Asia. The giraffe and the okapi (Giraffidae), two distinctive African species, are closely related to deer. The pronghorn (*Antilocapridae*), although sometimes called pronghorn antelope, is not a true antelope; it is the only survivor of a stock of ruminants that was very successful in the later part of the Tertiary Period in North America (about 2,500,000 to 65,000,000 years ago). The family Bovidae is primarily African and Eurasian, with a few members in North America. Bovids are advanced artiodactyls, many of which live in open grassland and semi-arid areas.

Importance to man. Artiodactyls have long been exploited by man for economic purposes. At Olduvai Gorge in East Africa there is clear evidence of the use of antelopes for food almost 2,000,000 years ago. In Europe during Paleolithic times (about 30,000 years ago) Cro-

Magnon man depended heavily on the reindeer. By this time the use of animals other than as food had become established; skins were used as clothing and footwear, and bones were used as tools, weapons, and accessories.

The domestication of animals was a major advance in human history. Domestication of herd animals probably arose gradually, perhaps before agriculture. Domesticated goats and sheep are first known from the Near East at some date close to 7000 BC. Cattle and pigs were domesticated at some subsequent date but certainly before 3000 BC. In South America the llama, now used for transport, and the alpaca, which provides a source of wool, were developed from guanacos by the Incas or their predecessors. The dromedary (*Camelus dromedarius*), domesticated in Arabia, was introduced into the Southwestern United States, southwestern Africa, and inland Australia in the 19th century. A large feral population now exists in Australia.

In addition to providing meat, milk, hides, and wool, artiodactyls have served man in a number of other ways. In Kashmir, the underfleece, or pashm, of the Siberian ibex (*Capra ibex*) and of local domesticated goats has been used as the basis for the manufacture of cashmere shawls. In southwestern France, pigs have been used to locate underground truffles (the fruiting bodies of certain edible fungi).

No group of mammals is more extensively hunted than the artiodactyls. Sport hunting of various deer supports a multimillion-dollar industry in North America and Europe. In many cultures hunting has been reserved for monarchs or the aristocracy. In the centuries after the Norman Conquest of England, the forest law provided severe punishment for the slaughter of deer and boars. Père David's deer (*Elaphurus davidianus*) of China now survives only because it was preserved first in the hunting park of the emperors of China and later by the Duke of Bedford after the slaughter of the Chinese herds at the end of the 19th century.

Wild ungulates were the primary source of meat for human populations long before the appearance of modern man. Prehistoric man hunted the large mammals of his environment with an ever increasing effectiveness that was certainly instrumental in his survival. The extent to which man was involved in the extinction of some of the larger Pleistocene animals (*i.e.*, those that were abundant 10,000 to 2,500,000 years ago) is still being investigated. There is now known to have been a wave of late Pleistocene extinction of large mammals, including artiodactyls; in North America this wave reached its zenith about 9000 BC. Many animals also became extinct in Africa, where long-horned buffalo and large relatives of hartebeests survived until very recently. More of the large mammals have survived in Africa than elsewhere, but the reason for their survival is not known. A second, probably final, wave of extermination of the larger mammals has taken place with the spread of European culture and firearms in the past 300 years. It has been marked by wanton slaughter and has ultimately produced an interest in conservation. It now seems, however, that the unprecedented demands on the environment being made by rapidly expanding human populations will result in a nearly complete extinction of large wild mammals.

NATURAL HISTORY

Behaviour. *Migration.* Many artiodactyls undertake seasonal migrations between their breeding grounds and feeding areas or between different feeding areas. They can then take advantage of the seasonal changes in different areas. This means that larger populations, and hence a larger biomass (*i.e.*, the total weight of all individuals in an area), can be supported than if all passed their lives in one area. The North American mule deer (*Odocoileus hemionus*) comes from its summer pastures at high altitudes as the first snow falls and returns at the end of winter, several weeks after the snow has melted.

Social behaviour. Although the popular image of artiodactyls is one of great herds numbering thousands of individuals, some species are solitary, and many others form only small family groups. The maternal family unit, in fact, is the most cohesive one, providing the basis for

Seasonal changes

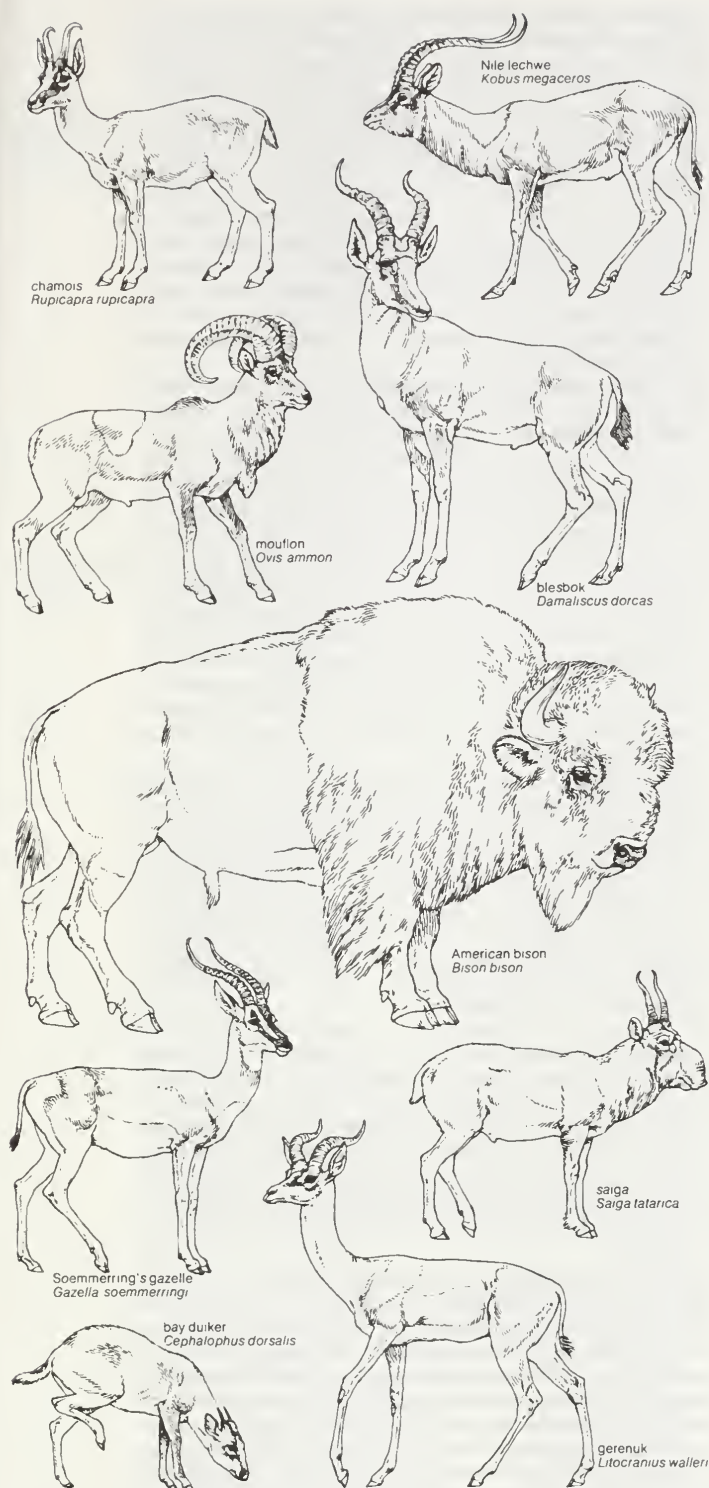


Figure 52: Body plans of bovid artiodactyls.

Drawing by R. Keane

The advantage of herd behaviour

herd formation. Most artiodactyls are more or less social, and grazing forms may be found in especially large aggregations. It appears that the practice of aggregating gives protection, favouring those members of the species that are the most active contributors to the gene pool (thus the most available to natural selection), since the individuals most frequently taken by predators are old, solitary males, males maintaining territories, and animals of either sex separated from the herd.

Social facilitation (the instigation of collective behaviour) takes place in herds. After one animal flees, all of the others flee, and the predator may thus not catch any. Social facilitation may also promote a restricted season for births; this helps survival of the young by denying these easy-prey individuals to predators through much of the year, and keeps the predator population lower than if young were

available throughout the year. Another advantage of herding is that the older generation in a herd can guide migrations to water, feeding areas, or mating grounds.

Females and young are usually in herds separate from those of the younger males, but territorial (the older, proven) males may accompany the females. There are some variations of this behaviour. In the Eurasian roe deer (*Capreolus capreolus*), for example, the basic unit includes the doe, her litter of two, and often the young of the previous year. During the rutting (mating) season males associate with females in heat but do not gather harems. The female herds of red deer (*Cervus elephas*) are separate from the males except in the breeding season, when the stag will defend his female herd against other males. Among cattle and related species, the males associate with the females and young, but the bulls are ranked below a so-called master bull, each defending its place within the rank order. Female hippopotamuses and their young form a group in water and have a favourite resting and basking sandbank. The males have their resting places around this area. Each male's rank in the social hierarchy determines how close to the females he may be.

There can be some flexibility of social organization within a species. During the rutting season the male Rocky Mountain goat (*Oreamnos americanus*) makes little effort to herd females within a fixed area if there is little snow, but he does drive off other males. When there is much snow, he neither fights other males nor defends individual females.

Forest-dwelling artiodactyls often live singly, as does the okapi (*Okapia johnstoni*) of central Africa; individuals meet only for mating. Female moose (*Alces alces*) with calves are intolerant of their own young of the previous year and of adults, so even small herds do not form.

The territory of an animal is an area from which the possessor attempts to exclude other individuals of the same species (and occasionally other species). An animal in an area lacking its own scent is more timid and ready to flee. Among solitary artiodactyls the territory holder defends an area sufficient to meet his needs for food and shelter. Among social artiodactyls the territorial system is interwoven with breeding activities, and territories are normally defended only by certain males. Other males are driven off, and a percentage of males are prevented from mating.

The most simple territorial organization among artiodactyls is that of the common wild pig (*Sus scrofa*), which lives within a home range including resting, feeding, drinking, and wallowing places. There is little sign of territorial defense, and the herd (called the sounder) may move to a new area. At the other extreme, male Uganda kob antelopes (*Kobus kob*) hold territories, for breeding only, that are as small as 15 to 30 metres (50 to 100 feet) in diameter. There are 30 to 40 territories on the breeding ground of a herd, and groups of females and young move about the territories despite the efforts of individual males to detain them. The semi-arid Serengeti plains of northern Tanzania contain nomadic aggregations of blue wildebeest (*Connochaetes taurinus*), males of which defend temporary territories only while an aggregation remains stationary.

In territorial defense an aggressive encounter between males is generally preceded by visual signalling of intentions. Chital deer (*Cervus axis*), for example, have several sorts of threatening displays. When sharp, potentially lethal horns appeared in early ruminants, intimidating displays rather than combats would doubtless have been favoured. Horns or antlers eventually functioned to maintain head contact during struggles rather than to bruise, slash, or gore. This stylized fighting, in which the competing males interlock horns or antlers and try to "outwrestle" each other, minimizes the danger of killing an opponent of the same species (conspecific). It evolved in two ways: further development of the wrestling, found in stags and some of the antelopes, and ramming, as in sheep. In sheep the horns are the sole organs of display. They increase in size throughout life and parallel the dominance order of the males, so that unnecessary fighting is minimized. Ramming may have intermediate forms; goats, for example, butt with a sideways hooking motion. In the fighting of hornless artiodactyls, such as pigs, the combatants may be

Territoriality

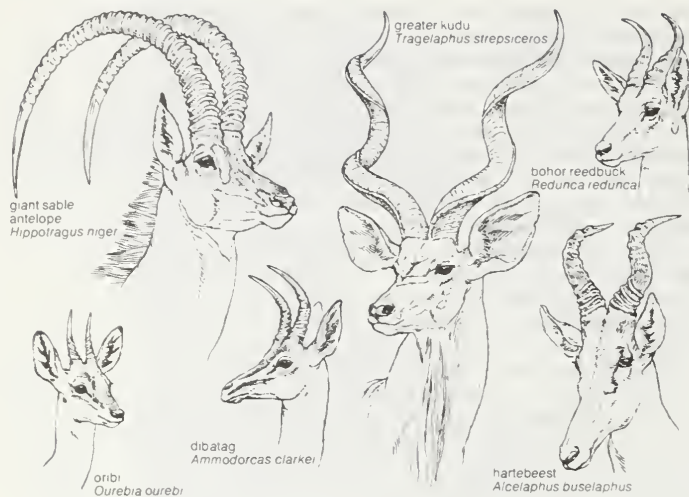


Figure 53: Horn variation in African antelopes.

Drawing by R. Keane

badly mauled or even killed. The fighting behaviour of camels retains primitive elements of biting, kicking, and neck wrestling.

Reproduction. Many advanced artiodactyls have elaborate courtship behaviour, a regular component of which is for the male to sniff or lick the female's urine, and afterward to raise his head slightly with upcurled lips. This behaviour, which has been called *flehmen*, apparently enables the male to recognize females in heat. In the mating ceremonies of tragelaphine antelopes (kudus, bushbucks, and others) the male follows the female, nuzzling her neck several times. When he mounts, he lays his neck along hers so that their heads touch. In Thomson's gazelle (*Gazella thomsoni*), following the *flehmen* behaviour, the male runs close behind the female and finally taps her hindleg with his foreleg. Similar leg contact also occurs in some other antelopes. Its function could be to test the female's readiness to mate, to habituate her to contact, or to heighten her readiness to mate. It appears to be equivalent to the neck contact of tragelaphines. During mounting, the male Thomson's gazelle holds his head high and does not touch the female's flanks with his forelegs; the pair may continue walking. This is probably a more advanced pattern of events than that in tragelaphines. The kob antelope has elaborate displays after mating. These and the specialized sexual displays seem to be a consequence of this species' tightly clustered territories on the mating grounds. Another pattern occurs in the normally solitary Indian hog deer (*Cervus porcinus*): as many as 20 or 30 aggregate loosely in a certain area, then females and males leave in pairs and usually remain together until they have mated. Mating in artiodactyls often intensifies toward dawn and dusk.

Gestation periods vary and are related in part to the size of the animal. They range from four months in the small chevrotain to 14 months in the Bactrian camel (*Camelus bactrianus*) and over 14 months in the giraffe. Females of normally gregarious species become solitary a few days before giving birth. The female chital, or axis deer, for example, remains near a patch of dense bush and high grass to which she can retreat if endangered. The female collared peccary (*Dicotyles tajacu*) withdraws to a burrow. The European wild pig gives birth in a rough nest.

In temperate regions, birth takes place in spring or early summer, and in tropical areas there are often more births during or just after the rainy season. The absence of a well-defined breeding season in a species may indicate less rigorous environmental conditions, which sometimes vary in different parts of a species' range. Warthogs have one restricted breeding season in most of eastern and southern Africa, while elsewhere two seasons or year-round breeding have been recorded. The breeding season of the waterbuck (*Kobus ellipsiprymnus*) is continuous in Uganda, but in Zambia its breeding season shows a sharp peak at the height of the rains.

Most modern artiodactyls have one young at each birth,

but there are some well-known exceptions among ruminants. The Chinese water deer (*Hydropotes inermis*) bears twins or triplets, but during gestation carries even more fetuses; early records (now known to be incorrect) of large litters were based on observations of dead pregnant females containing the large number of fetuses. The mule deer, white-tailed deer (*Odocoileus virginianus*), roe deer, pronghorn (*Antilocapra americana*), nilgai (*Boselaphus tragocamelus*), four-horned antelope (*Tetracerus quadricornis*), and saiga (*Saiga tatarica*) commonly bear twins. In the white-tailed and mule deer and in the saiga, a higher percentage of twins are borne by the older females; this is probably true in other species. The number of young is usually three in the warthog, five in the European wild pig, and two in peccaries.

The female wild pig almost ignores her young, which free themselves from their birth membranes and seek a teat. Female camels show comparatively little maternal attention and do not eat the afterbirth (the fetal membranes and placenta). Ruminants generally eat the afterbirth, as well as the dung and urine of the young, thus helping to prevent discovery of the young by predators. Licking of the young tends to facilitate its recognition by the mother. An artiodactyl is normally precocious (well developed) at birth and may weigh one-tenth as much as its mother. An extreme example of precocity is the wildebeest calf, which rises within five minutes of birth, follows its mother within another five minutes, and can move as fast as an adult in 24 hours. Young deer fawns "freeze" during danger but rejoin the herd when the danger is long past or when retrieved by the mother.

Pigs and hippopotamuses are weaned after a few months, but among higher artiodactyls, lactation lasts longer. Wildebeest, for example, suckle for almost a year, although they start to eat grass when only a few days old. This may either maintain a bond between parent and offspring and form the base for larger social groupings or help to "develop" the four-chambered stomach. Higher artiodactyls eat soil when they begin to eat solid food, probably to establish a normal flora and fauna in the rumen (the first of the four stomach chambers).

Locomotion. Artiodactyls are preyed upon by carnivores and therefore need speed and agility to escape death. They have an added disadvantage in the sheer weight of their very large stomachs, which they need in order to digest plant food. Running ability reaches an extreme in advanced artiodactyls living in open country. The hippopotamus, with an adult weight of 2,500 to 3,000 kilograms (5,500 to 6,600 pounds), is the only living artiodactyl big enough to need heavy, pillar-like limbs for support.

In the normal walking of artiodactyls the legs move in the following order: (a) left front, (b) right rear, (c) right front, (d) left rear. This basic pattern is masked in faster walking or trotting by each foot being lifted off the ground before the one ahead of it in the sequence reaches the ground, resulting in telescoping the first (a and b) and second (c and d) pairs of movements. In galloping or fast running the two front legs leave the ground one immediately after the other, then the two back legs. The chief propulsive force in locomotion comes from the back legs, except in the giraffe (*Giraffa camelopardalis*), in which the front legs provide the main propulsive power.

Camels often amble, both legs of each side moving together, and the giraffe and the okapi always use this walking gait. Here the middle two (b and c) and the first and last (a and d) actions of the normal walking pattern occur together. The giraffe, having a short body and great height, could not adopt the normal ruminant gait without tripping. The long neck moves back and forth in time with the strides and helps smooth the movement. Galloping by the giraffe is of the normal ungulate type.

Artiodactyls living among bush or rocky cover may develop a bounding sort of gait in which the legs are pulled up very sharply during each stride. Deer and some antelopes are examples. When walking, species in such habitats are supported by the diagonally opposite legs for a greater length of time in each stride than are fast-running, open-country ruminants. This is a more primitive stable

Precocity
in offspring

Gestation
periods

position and allows an easier leap from hidden danger. Some bovids, notably goats in Eurasia and the klipspringer (*Oreotragus oreotragus*) of Africa, are especially agile on rocky slopes and precipitous ground.

The maximum speeds of some artiodactyls are: warthog, 48 kilometres (30 miles) per hour; camel, 14–16 kmph (9–10 mph); giraffe, a little over 48 kmph (30 mph); Cape buffalo (*Syncerus caffer*), 56 kmph (35 mph); Thomson's gazelle, 80 kmph (50 mph).

Ecology. *Food habits.* Most artiodactyls are closely tied to the resources of their environment. They are dependent, for example, on feeding areas not being covered by too much snow or shrivelled under a drought, and on the regulating effects of fire or other herbivores on the seasonal succession of vegetation. Various grazing species feed on grass at different heights. Browsers, those that feed on the foliage of shrubs and trees, show more extreme variation in feeding height, the maximum being that of the giraffe.

Herbivorous animals need less initiative and intelligence to collect food than do the meat-eating, hunting carnivores, but digestion is more difficult. Advanced artiodactyls have evolved the ability to bolt food and to ruminate it (chew it more thoroughly) at a later time or while resting in an area where they may be less obvious to predators and can conserve energy. Tropical artiodactyls frequently have adaptations for water conservation, having developed to a high degree internal physiological regulation (homeostasis).

Primitive artiodactyls were probably omnivorous but favoured plant foods, a characteristic still found in pigs. The latter dig with the snout and, to a lesser extent, with the front legs and upper tusks (canine teeth). The wart-hog of Africa (*Phacochoerus aethiopicus*) has a modified method of gathering food. When food is scarce it forages for young grass shoots under very low bushes; its tusks and localized thickening on its skin protect the eyes and muscles from thorn damage, and small incisors enable it to pluck food.

Hippopotamuses (*Hippopotamus amphibius*), although they spend a great deal of time submerged in lakes or rivers, do not feed in the water. They graze at night, wandering over well-used trails, sometimes far from water, often damaging crops.

Most members of the camel family are found in arid habitats. The vicuña (*Lama vicugna*) of the South American Andes lives at high altitudes where it grazes on soft grasses and herbs. It has much the same food requirements as domestic sheep.

Chevrotains live in dense undergrowth close to water or in marshes, where they browse on soft vegetation, roots, and tubers, following a way of life probably not unlike that of their ancestors.

The other ruminants browse or graze. They may take many plant species in the course of the year, but at any one season a large part of the diet consists of only five or six plants. Some ruminants are strongly specialized. The reindeer of the Arctic (*Rangifer tarandus*), for example, eats a variety of sedges, grasses, and herbaceous plants in summer but, as the long winter approaches, gradually shifts to a diet of lichens. It uses its front feet to scrape snow away from lichens to a depth of about 60 centimetres (two feet). The females are unique among deer in possessing antlers, which are thought to help them get scarce food in late winter by driving off the males that have by then shed their antlers. Reindeer may eat lemmings. The red deer, on the other hand, has catholic feeding habits. In woods it browses on lichens, berries, fungi, and the leaves of most deciduous trees; in open country it eats grass, heather, berries, and lichens. Shrubs and trees are used more in winter. When the red deer lives in the same areas as other ruminants it can be a serious competitor for food.

Grasses form a substantial part of the diet of many ruminants. Young grass consists of about 5 percent protein, 1 percent fat, 3 percent minerals, and 20 percent carbohydrates; the remaining percentage is water. The most noticeable changes as grass ages are an increase in carbohydrate content to 75 percent and a large decrease in the amount of water. Such food, especially when coated with silica, as are many grasses, or when covered with dust, would be impossible for nearly all nonruminant herbivores to eat

or digest. The major evolutionary trend in ruminants has been to make use of grasses and grasslands, and the higher ruminants have evolved largely in adaptive balance with one another. This adaptive balance was shown during a study of the change from plains to thickets of scrub growth in an area in the eastern Congo over a period of about ten years. There was an accompanying decrease in numbers of antelopes and warthogs, no change in buffalo, and an increase in elephants and hippopotamuses.

There is not usually a one-to-one dependence of any artiodactyl species on one plant. The plant species that constitute the major part of the diet may vary with the season, and similar parts of different plants may be eaten in preference to other parts of the same plant. Food resources in an area are thus parcelled out among the various artiodactyls present. Sometimes behavioral differences minimize competition between closely related species in the same area. A study has shown that in central Africa the roan antelope (*Hippotragus equinus*), a grazer, favours open areas with taller, ranker perennial grasses and is more or less sedentary within a small area; the sable antelope (*H. niger*), also a grazer, prefers savanna woodland or the edges of open areas, and herds follow a more or less cyclic annual route over an area of about 500 square kilometres (200 square miles). When pasturage is restricted, sheep will cut grass very short, and goats will damage trees and bushes. An American zoologist, George B. Schaller, has observed that, in Kanha Park in central India in the hot season, blackbuck (*Antelope cervicapra*) continue to graze on grass shoots in open areas; chital deer seek out tender grass blades, especially along forest edges, and also feed on leaves and fruits; barasingha (*Cervus duvauceli*) eat dry and moderately coarse grass along ravines; sambar deer (*Cervus unicolor*) browse on leaves and crop coarse grasses in the forest; and gaur (*Bos gaurus*) graze on tall, coarse grass and break down saplings to get at the leaves. The choice of habitat also varies: chital avoid steep terrain and forests with an unbroken canopy; blackbuck require less water than the others and thus remain in drier regions; sambar and gaur are less specialized in habitat requirements, and both are active primarily at night; barasingha prefer reed beds but also enter forests and climb hills.

It has also become evident that grazing successions are one of the mechanisms that enable the maximum use to be made of environmental resources. On the Serengeti plains, for example, the wildebeest grazes on ground already covered by the zebra and leaves the grazed grass in a condition suitable for the Thomson's gazelle. Interactions take place between artiodactyls and some plant species. It has been noted in the Tarangire area of northern Tanzania that *Acacia* seedlings germinate only where the impala (*Aepyceros melampus*) has left its dung. In parts of southern Peru plants growing on or close to the dung of the vicuña are different from those of the surrounding pasture.

Artiodactyls often favour the boundary zone between habitats. In Rhodesia, Lichtenstein's hartebeest (*Alcelaphus lichtensteini*) is usually found at the edge of clearings adjacent to woodland.

Areas of distribution. Some artiodactyls have surprisingly small ranges; Hunter's hartebeest (*Beatragus hunteri*) and the dibatag (*Ammodorcas clarkei*), for example, are found in two very restricted areas in eastern Africa. Others have extremely large ranges, such as the roe deer, which lives from the western shores of Europe to the eastern shores of Asia, or the red deer, which is found in a similar band across Eurasia and is regarded by many as conspecific with the North American wapiti or elk (otherwise called *Cervus canadensis*). Sometimes a considerable area may be occupied by a chain of related species, an example being the oryxes; the beisa and gmsbok (races of *Oryx gazella*) occur in South and East Africa, the scimitar-horned oryx (*O. dammah*) in West Africa, and the Arabian oryx (*O. leucoryx*) in Arabia.

It is well known that climate is one of the factors limiting the ranges of artiodactyls. A number of South African antelopes differ, at the species level, from their ecological counterparts farther north in Africa. The bontebok and blesbok, races of *Damaliscus dorcas*, are found in the south and the sassaby (*D. lunatus*) farther north; the black

The importance of grasses in the ruminant diet

The influence of climate on the distribution of artiodactyls

wildebeest (*Connochaetes gnou*) occurs in the south and the blue wildebeest (*C. taurinus*) farther north. This probably is a result of climatic or climatically influenced factors; each species evidently functions best in a certain temperature and aridity range. Wide distributions can occur more easily along lines of latitude than they can by spanning the tropics to temperate or polar regions. Species that cross lines of latitude are often associated with mountain chains, examples being the Rocky Mountain goat, with its wide latitudinal range in western North America, and the goral (*Nemorhaedus goral*), found from Indochina to the Amur River. Climatic effects on distributions sometimes occur with regard to altitude. In Central Asia, the goat (*Gazella picticaudata*) is found in valleys from 3,000 to 3,660 metres (10,000 to 12,000 feet) above sea level, the chiru (*Pantholops hodgsoni*) and the yak (*Bos mutus*) are on the very high steppe between 5,500 and 6,100 metres (18,000 and 20,000 feet).

South America has a more impoverished artiodactyl fauna than Africa, being limited to deer and camelids. This arises in part from the late arrival of the artiodactyls (deer in middle to late Pliocene, about 4,000,000 years ago, camelids perhaps a little later) and in part because a number of large rodents compensate for the shortage of large herbivores. The cervids in South America have not shown the same capacity for radiation in open country as have bovids in the Old World.

Influence
of pests on
distribution

The areas of distribution and numbers of individuals are determined by complicated interweaving of effects not yet completely understood. Bloodsucking flies are thought to be the main reason that red deer in Scotland ascend to higher feeding grounds in June, and reindeer are afflicted by horse flies (*Tabanus*) and other dipteran pests. It is questionable whether the level of artiodactyl populations is controlled by predation, by availability of food, by reproductive rate, by disease, by climate, or by competition, insofar as these can be regarded as separate factors. It is known that undernourishment increases the susceptibility of an animal to the effects of parasites. If such an infected animal, say a pig, is caught by a leopard, it would be an oversimplification to assign a single reason for its death; it could have died from starvation, parasites, or predation. There is no evidence that artiodactyls are affected more than marginally by predators during most of their mature lives. Mortality is greatest among juvenile and aged animals. In a study of central African warthogs, it was estimated that a 60 percent loss occurred during the first six months of life in an expanding population and 95 percent in a declining one. Although predation was thought to be the main cause, another was the fact that the piglets had only limited control over their body temperatures and were thus more at the mercy of environmental temperature change. Food supply may sometimes be decisive, either directly or through the indirect action of intermediate agencies such as drought. The year 1961 lacked long rains, causing a severe shortage of forage in the Nairobi Game Park in Kenya. Many antelopes died of starvation, populations fell, and those of the blue wildebeest had not recovered nine years later, perhaps for reasons unconnected with the initial drought. Disease has generally been considered to have only a secondary importance in regulating numbers.

Thickness of the snow cover in winter is a very important factor for Asian artiodactyls. The saiga, for example, cannot move in snow deeper than about 40 centimetres (16 inches), and the wild sheep *Ovis ammon* in snow deeper than 60 centimetres (24 inches), at the most. The snow may have other effects; a layer of ice on top of snow may damage an animal's legs and weaken the animal to the extent that it is caught by a predator. Saiga may be unable to dig through even a shallow layer of compacted snow. Hoarfrost on vegetation is especially dangerous when prolonged or when it occurs in consecutive winters, though elk may escape the worst effects by feeding in winter on bark and high shoots. Massive periodic mortalities among Palearctic (Eurasian) ungulates in winter have been known since ancient times. The saiga has adapted to these crises by migrating great distances in a short time away from snowstorms or from areas where fodder is short. It also has

a very rapid maturation to a reproductive state, ensuring that populations will build up after heavy mortalities.

Population density over the range of a species is affected by social behaviour, such as the effects of territoriality, dispersal of the young, and whether the species lives in herds. Fecundity may be reduced in overcrowded conditions by effects on reproductive control mechanisms, reduced viability of the young, or retarded maturation.

FORM AND FUNCTION

General structure. Artiodactyls have larger stomachs and longer intestines than carnivorous animals because plant food is less easily digested than meat. The necessity of escaping predators and the handicap of a heavy digestive system have resulted in limb bone adaptations.

In all artiodactyls the main weight-bearing axis of the leg passes through the third and fourth toes together. This has been called paraxonic support and is contrasted with the mesaxonic limb support of the other great order of herbivorous mammals, the perissodactyls (rhinoceros, horse, tapir), in which the weight-bearing axis passes through the third or central toe alone. As artiodactyls evolved there was increasing development of the third and fourth toes and a parallel decline of the second and fifth toes flanking them. Progressive simplification of limb extremities has characterized the evolution of the artiodactyls, and even in the earliest known artiodactyls, the pollex and hallux (corresponding to the big toe and thumb of man) were already rare.

Comparison of
odd- and
even-toed
ungulates

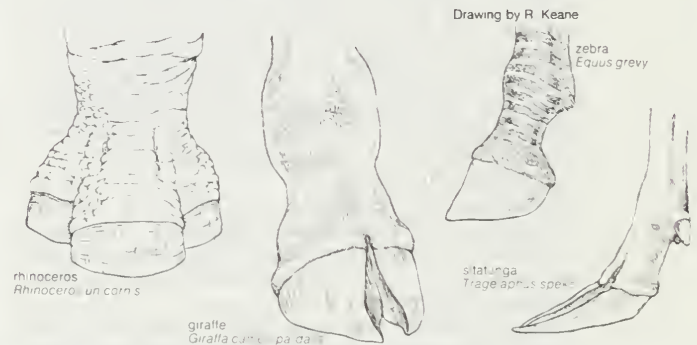


Figure 54: Comparison of even-toed (artiodactyl: giraffe and sitatunga) and odd-toed (perissodactyl: rhinoceros and zebra) ungulate feet.

The other main morphological characteristic of artiodactyls is that the astragalus, one of the bones in the ankle, has upper and lower rounded articulations (areas of contact of bones) and no constricted neck, instead of simply one rounded articulation above a neck, as in other mammals. This character is so basic to artiodactyls that it has not developed very much within the known history of the order, having already been present in long extinct members. The artiodactyl astragalus also has an articulation on its rear surface for the calcaneum (heel bone). The three articulations are in nearly parallel planes, allowing the astragalus to rotate vertically.

Other features of the limbs, skull, and dentition distinguish artiodactyls. The ulna (posterior forearm bone) and fibula (posterior bone of the lower leg) have become reduced. The humerus, the upper bone of the foreleg, is large and has a large protrusion, the greater trochanter, to which muscles are attached. The femur, the upper bone of the hindleg, has a large greater trochanter and a second, lesser trochanter, but lacks the third trochanter characteristic of perissodactyls. There are typically 19 thoracic and lumbar (upper and lower back) vertebrae. The separate lumbar region of the spine is retained with its forwardly directed transverse processes (lateral projections on the vertebrae). There is no clavicle, or collarbone, in the shoulder girdle. The hip girdle shows fore-and-aft elongation and a well-developed ischium (upper anterior bone of the pelvis). There is never a penis bone.

The large tongue is very mobile and can be thrust forward. The brain is moderately developed, with folding of the surface of the cerebral hemispheres variably developed, often less in small artiodactyls than in large ones.

The olfactory region of the brain is well developed and hearing is acute. The brains of earlier artiodactyls, such as the extinct entelodonts, were smaller than those of later forms. There are often scent glands on the head and body.

Specializations of the head. The skulls of pigs and peccaries lack a complete bony bar behind the eye (postorbital bar) as in most suiform artiodactyls and the early camels. The hippopotamuses, most camels, all ruminants, and two fossil suiform groups (entelodonts and oreodonts) have a complete postorbital bar. Any surface exposure of the petriotic bone (bone around the ear) on the skull is called the mastoid, and skulls without such a surface exposure are described as being amastoid. Amastoid skulls are found in most suiform groups (including entelodonts, anthracotheres, and all living suiform groups); mastoid skulls occur in some early suiform groups, oreodonts, and all remaining artiodactyls that have lived since the end of the Eocene Epoch (about 38,000,000 years ago). Hippopotamuses have many modifications for aquatic life—large lungs, eyes and nostrils on top of the head, nostrils that can be closed by muscular control, and small ears. They are able to remain submerged for at least five minutes.

Horns and antlers Pigs, peccaries, hippopotamuses, camels, and chevrotains have no horns or antlers. In the early Miocene, Old World ruminants related to giraffes and deer first developed such appendages. The majority of deer have antlers, defined as solid, bony, branched outgrowths of the frontal bones, present only in the males (but also in female reindeer) and shed seasonally. They are not covered by a horny sheath but, during a growth period of about four months, have a fine-haired skin or "velvet." The antlers have two basic branches, the anterior or brow tine, and the posterior branch or beam. The brow tine is unbranched, except in Père David's deer, in which both it and the beam are branched, the brow tine forming the dominant part of the antler. Antlers are specialized sex characters used for fighting by males in the rutting season and to scrape or slash at trees and bushes for territorial marking.

A study of the chital deer showed that antlers increase in size up to the seventh year, remain at a constant size until the ninth year, then decline. The horn of bovids consists of a hollow, unbranched horny sheath (formed of modified skin like fingernails and toenails) that fits over a bony core; horns are often present on both sexes. If such a horn is accidentally lost it is not regenerated; this is unlike the situation in deer, in which normal shedding is followed by regrowth. In the giraffe, but not in the okapi, horn growth is mainly from the parietal bone. The pronghorn has horns in both sexes. The sheaths are shed each year after the breeding season, and new ones develop under the old ones. The sheath is two pronged, but the underlying bony core is unbranched.

Teeth. There is a complete set of teeth in early artiodactyls and in modern pigs of the genus *Sus*, consisting on each side of three upper and lower incisors, an upper and lower canine, four upper and lower premolars, and three upper and lower molars. There has been a tendency toward reduction of the front teeth and development of a gap (diastema) between them and the back teeth. There has been very little tendency for the premolars to molarize, and the first premolar often disappears. Early forms had five-cusped upper molars, but the fifth cusp (protoconule) disappeared early.

Members of the suborder Suiformes have the full complement of incisors and canines, except for peccaries, which lack the lateral pair of upper incisors. Hippopotamuses have continuously growing incisors and canines, the lower canines being very large.

The canines of pigs grow continuously. In this group the canines are weapons for offense and defense, the sharp cutting edges of the lower canines being maintained by wear against the uppers. Young camels retain the full complement of front teeth, with three incisors and one canine in the upper and lower jaws; the upper incisors are extremely small. In the upper jaw of the adult only the rear incisor and canine are present. The vicuña has continuously growing lower incisors.

The molars of pigs are low crowned (except those of

the warthog) and have many cusps; those of peccaries are more simple. Peccaries have one less premolar than pigs; camels also have reduced premolars. Chevrotains have rather flattened lower premolars but have incipiently selenodont molars; *i.e.*, in which the cusps are drawn out into longitudinal crescents. Premolars of ruminants are wider, and the molars definitely selenodont. In many bovids and the pronghorn, but not in giraffes or deer, the molars are markedly high crowned.

Limb adaptations for fast running. Adaptations for fast running reach an extreme in advanced artiodactyls living in open country. In addition to the increased rotation of the astragalus, which increases the propulsive thrust at the ankle and enables a quicker recovery at the end of a stride before starting the next one, there are other features that help to increase the speed of striding. The legs of most camels and ruminants have lengthened, especially in the lower parts; the number of toes, or digits, in the feet is reduced from the original mammalian five, and ruminants walk on the tips of their toes. The muscles are inserted high on the legs; only tendons pass lower, so that a large mass is not concentrated near the tip of the limb, where its inertia would restrict speed of movement. Muscle contraction is fast. The movement of each leg is almost limited to a fore-and-aft plane. Emphasis on the fore-and-aft articulations between the limb bones is especially pronounced in many bovids, the alternating bones in the wrist (carpus) and ankle (tarsus) taking the strain of impact on uneven ground.

Pigs have four toes on each foot, but only two of them touch the ground. Their limbs are short and not very advanced. Peccaries have lost the outer accessory hind hoof in the back leg. All four toes of each foot of hippopotamuses touch the ground, and the terminal phalanges have nail-like hoofs. The toe bones of camels are completely enclosed in hardened, horny hoofs, and lateral toes spread across the broad pad which aids in walking on desert sands. Chevrotains have four hoofed toes on each foot; deer often retain the first and second phalanges (sections) of their lateral toes; but all bovids have lost the bones of their lateral toes.

The fibula bone in the back leg and the ulna in the front leg have been reduced in different artiodactyl lineages. Both are still complete in pigs and hippopotamuses, although the fibula is slender. In most other artiodactyls, the lower end of the fibula has survived, and the upper end is occasionally found, but always less noticeably. In camels the ulna has fused with the radius. Pigs, hippopotamuses, and camels have separate navicular and cuboid bones in the ankle, and magnum and trapezoid bones in the wrist; other artiodactyls have a fused naviculo-cuboid and magnum-trapezoid. In chevrotains and some deer, the adjacent ectocuneiform is sometimes joined with the naviculo-cuboid.

The artiodactyl method of limb support through the third and fourth toes, with the attendant lengthening of lower limb bones, has frequently led to a fusion of the two principal metacarpal and metatarsal (midfoot) bones in the forelegs and hindlegs, respectively, forming cannon bones. The nearest approach to a cannon bone in the living Suiformes is the proximal fusion (*i.e.*, at the upper ends) of the two central metatarsals in peccaries. Camels have front and rear cannon bones, but the fusion does not extend right to the bottom, the lower articular surfaces being less pulley-like than in ruminants. There is a hind cannon bone in all chevrotains and, in addition, a front one in Asiatic species (*Tragulus*). All other living artiodactyls have front and rear cannon bones. Lateral metatarsals and metacarpals survive in chevrotains; splints of lateral metacarpals often survive in bovids; and either upper or lower splints of metacarpals in deer.

Modifications of the skin. *Hair and coloration.* Pigs are covered with rather sparse, coarse hairs, and peccaries with a denser coat of coarse hairs. Except for those of the warthog and the babirusa (*Babyrusa babirusa*), piglets have longitudinal stripes or flecks. Hippopotamuses are naked. Tragulids have light-coloured flecks and stripes in their fur. The coats of camelids and deer are much thicker in species living toward the polar regions, at great heights,

The distinction between horns and antlers

Modifications of the foot

or in deserts, but are not noted for striking colours or patterns. Many young deer and the adults of a few species have pale flecks and stripes, and some South American deer have reddish fur. Antelopes have a wider range of coat colours, and some are strikingly marked; e.g., the oryxes, bontebok, and blesbok of southern Africa.

Scent glands. External glands occur in various places on artiodactyls. Preorbital glands, immediately in front of the eyes, are present in the giant forest hog (*Hylochoerus meinertzhageni*), in all cervids except the roe deer, and, among the bovids, in duikers, many neotragines, gazelles and their allies, and the hartebeest group. These glands are apparently required in small forest forms and have disappeared in many, but not all, open-country forms. In some, the glands are definitely connected with territorial marking; a firm object is marked by rubbing, soft vegetation by swinging the head gently from side to side. Foot, or pedal, glands are present in the African bush pig (*Potamochoerus porcus*), camels, tragulids, the pronghorn, some bovids, and on the back legs only of most American deer.

Inguinal (belly) glands are found in bovids, there being two in sheep, saiga, chiru, gazelles, duikers, and blackbuck, and four in members of the tribes Reduncini and Tragelaphini. Carpal (wrist) glands are present in some pigs, some gazelles and allies, and the oribi (*Ourebia ourebi*). Glands in other positions are rather less frequent, but postcornual ones (behind the horns) occur in the Rocky Mountain goat, the pronghorn, and the chamois (*Rupicapra rupicapra*), supraorbital ones in muntjacs (several species of *Muntiacus*). There are jaw glands in the pronghorn; neck glands in camels; dorsal glands on the back of peccaries, pronghorn, and springbok; and preputial glands (in front of the genital region) in several pigs, grysbok (*Raphicerus melanotis*), and the musk deer. Tail glands are found in musk deer, pronghorn, and goats; tarsal glands in pronghorn and American deer; and metatarsal glands in camels, some deer, and the impala. Pronghorn, blackbuck, gazelles, and oribi are thus particularly well equipped with glands. The use of such glands, apart from the use of preorbital glands in some species for territorial marking, is a matter for conjecture. Chital deer, when alarmed, thump the ground several times with their hind feet, which possess glands; the scent remaining on the ground may function as a danger signal. In general, mammals often mark with their glands when they are threatening other individuals of their own species.

Digestive system. The higher artiodactyls feed only on plant matter, which consists largely of cellulose and other carbohydrates and water. This necessitates adaptations of the structure and functioning of the stomach and intestines. Even pigs have enlarged stomachs—they have a pouch near the cardiac orifice (the upper opening) of the stomach—and in peccaries the stomach is more complicated. In hippopotamuses the stomach is divided into four compartments, and micro-organisms ferment food as part of the digestive process. Unlike pigs, hippopotamuses have lost the cecum (a blind pouch) further on in the gut.

In the most advanced ruminants, the much enlarged stomach consists of four parts. These include the large rumen (or paunch), the reticulum, the omasum (psalterium or manplies)—which are all believed to be derived from the esophagus—and the abomasum (or reed), which corresponds to the stomach of other mammals. The omasum is almost absent in chevrotains. Camels have a three-chambered stomach, lacking the separation of omasum and abomasum; the rumen and reticulum are equipped with glandular pockets separated by muscular walls having sphincters (valves) and glands. The esophagus opens into the rumen, not into the area between rumen and reticulum; these and other differences suggest that camels evolved the ruminating habit independently of the true ruminants. The total stomach of the domestic ox (*Bos taurus*) occupies nearly three-quarters of the abdominal cavity, and, even in medium-sized cattle, the rumen alone can have a capacity of 95 to 285 litres (25 to 75 gallons), having undergone a tremendous growth in early life, with the changeover from a milk diet.

Food taken into the rumen is later regurgitated into the mouth and completely masticated, then swallowed again

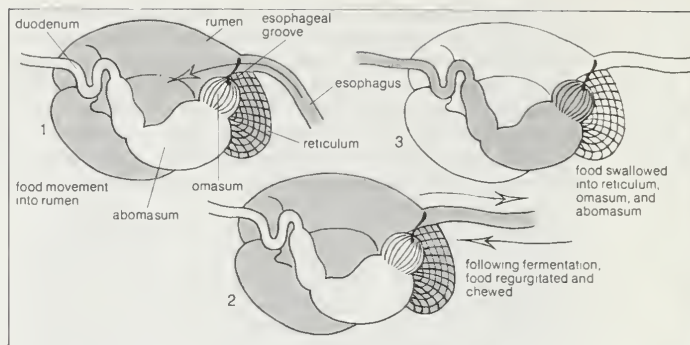


Figure 55: Stages of digestion in the four-part stomach of a representative ruminant.

Drawing by R. Keane

and passed to the reticulum, omasum, and abomasum. The regurgitation and chewing in the mouth is called rumination.

In the rumen many different species of minute protozoans (animals) and bacteria live without free oxygen. The digestion of the cellulose of plant cell walls is the main function of the fauna and flora in the rumen, since mammalian digestive juices are incapable of digesting cellulose. The contents of the plant cells are thus released for digestion. Large volumes of saliva are secreted into the rumen to help digestion. Soluble products of microbial action, mainly fatty acids, are absorbed through the rumen wall. In the omasum, some fatty acids and 60–70 percent of the water are absorbed; in the abomasum gastric juice containing hydrochloric acid is secreted, as in an ordinary mammalian stomach.

In the rumen any ingested protein is degraded into fatty acids and ammonia; the ammonia and other simple nitrogen-containing substances are used by the micro-organisms for their own cell-protein synthesis. These organisms are ultimately digested in the abomasum and small intestine, thus providing the ruminant with protein.

Many artiodactyls are adapted to living in conditions of water shortage. The best known and one of the most spectacular examples of this is the camel. Its body temperature can fluctuate according to the outside temperature, thus minimizing water loss through sweating; it excretes rather dry dung and a concentrated urine (*i.e.*, high in urea and low in water) and is not seriously weakened by as much as a 25 percent dehydration in its body, since water is not withdrawn from the bloodstream and the continuing circulation avoids any buildup of excessive internal temperatures. The thick coat hinders the inward transference of heat from the environment (the temperature of which may often exceed the animal's body temperature); a thirsty camel can take in water very rapidly. Oryxes and gazelles are antelopes noted for needing little water, the dorcas gazelle (*Gazella dorcas*) in the Sudan depending on leaves of *Acacia* bushes for its water. The zebu (a form of domesticated cattle) needs less water than most temperate climate breeds.

Reproductive specializations. The testes of male artiodactyls descend outside the body cavity but may regress into the abdomen in the nonbreeding season. Female pigs have many teats, but ruminants have only two to four (although domestic cattle occasionally have as many as six). Among the bovids, the alcelaphines (hartebeests, wildebeests, and relatives), gazelles, and some caprines (sheep, goats, and relatives) have two, the rest have four.

The unborn mammal within its mother breathes, feeds, and excretes through an organ called the placenta, which is connected with the tissues of the mother's uterus (womb) wall. Hippopotamuses and pigs have an epitheliochorial placenta, a layer of fetal tissue merely pressed close against the uterus wall, but camels and ruminants possess a syndesmochorial placenta, in which the epithelium of the maternal tissues is eroded to facilitate intercommunication. This is an advance over the epitheliochorial placenta, but the artiodactyls are not particularly advanced, when compared with other mammals, in which there may be still closer association of maternal and fetal blood vessels

The epitheliochorial and syndesmochorial placentas

Ruminant stomach

(endothelial and hemochorial placentas). Even in many syndesmochorial placentas the uterus lining may be wholly or partly restored before the end of pregnancy. Although there is no erosion of maternal tissues in the epitheliochorial placenta, the capillaries beneath the fetal and maternal surface layers may pass just beneath the surface layers, making them thin. The actual fingerlike processes (villi), through which the placenta contacts the uterus, are evenly distributed ("diffuse" placentas) in hippopotamuses, pigs, camels, and tragulids; in higher artiodactyls they are in pockets or groups called cotyledons ("cotyledonary" placentas). It is interesting that there are few of these cotyledons in deer—for instance only five in Père David's deer—but many in giraffes and bovids (up to 160 or 180 in giraffes and goats). The musk deer (*Moschus moschiferus*) is exceptional among deer in retaining a diffuse placenta.

EVOLUTION AND PALEONTOLOGY

The artiodactyls can be traced back to a probable descent from a group of early generalized mammals called condylarths, and were certainly distinct by the Eocene Epoch, which ended about 38,000,000 years ago. Fossil artiodactyls can be more or less convincingly classified in three suborders; the more primitive Suiformes, centred around pigs, the Tylopoda, centred on camels, and the Ruminantia or ruminants. The most primitive artiodactyls are the suiform group Palaeodonta, which had four functional toes on each foot, primitive, low-cusped cheek teeth, and the typical artiodactyl astragalus. The artiodactyls became more prominent in the Oligocene (between about 38,000,000 and 26,000,000 years ago) with a decline of the then dominant perissodactyls, and the later history of artiodactyls appears as successive waves of groups, each better adapted than its predecessors to the changing environment. In the suiform line, the earlier palaeodonts are succeeded by other groups such as the entelodonts, giant "pigs" of the European and North American Oligocene, characterized by very large skulls (some nearly a metre [three feet] long), very small brains, and a large, bony flange below the eyes. The functionally two-toed ruminants succeeded four-toed suiforms in the Miocene, and within the Old World ruminants of the bovid subfamily Caprinae, the zenith of the tribe Caprini, for example, followed that of the mainly Pliocene tribe Ovibovini.

The artiodactyls had an interesting history in North America through the Tertiary Period. Some forms, such as the entelodonts, were shared with the Old World, but others were characteristic of North America. One very prominent New World family was the merycoidodonts (or oreodonts), which lasted until the early Pliocene (about 6,000,000 years ago). They had somewhat piglike proportions, short faces, a large upper canine and a caniniform first lower premolar, and selenodont molars. A close relative, *Agrichoerus*, had clawed feet, the function of which remains uncertain.

Camelids evolved in North America and, at or toward the end of the Tertiary, spread into South America and into the Old World. By the end of the Pleistocene they all became extinct in their homeland, just as horses did. The hypertragulids were a mainly Oligocene group of chevrotain-like forms related to the Protoceratidae. The latter had horns above their noses, a position unique among artiodactyls, as well as in the usual position. The North American Miocene (26,000,000 to 7,000,000 years ago) produced some ruminants, such as *Blastomeryx*, that are hard to distinguish from the early palaeomerycine relatives of giraffes and deer in the Old World, which, with the North American groups, constitute the family Palaeomerycidae. Some developed horns, and the dromomerycine *Cranioceras* even had a third horn above the back of its skull. During the Miocene and Pliocene there finally appeared relatives of the surviving pronghorn, an example being *Merycodus*. Many of these North American groups have parallels with Old World groups, and the subject of North American artiodactyl evolution is of great interest. Only further finds will indicate whether *Blastomeryx*, the dromomerycines, *Merycodus*, and the pronghorns evolved from hypertragulids already in North

America or sprang from some immigrant ruminant and, if the latter, whether the supposed hypertragulid *Leptomeryx* could be such an immigrant ruminant. It is uncertain whether the hypertragulids are nearer the tragulines or the camels, and how close the oreodonts are to the anthracotheres. Of the great New World radiation there survived after the Pleistocene only three or four camelid species and the pronghorn (deer and bovids in the Americas are immigrants), whereas in the Old World as little as 200 years ago, Eurasia and Africa had abundant deer and antelopes.

Until the Miocene there were some archaic artiodactyls in Europe, the xiphodonts, which have cautiously been taken as tylopods, and the cainotheres and anoplotheres, which are classified near anthracotheres.

A possible ruminant ancestor was *Archaeomeryx* from the upper Eocene of China, a small animal that already had a fused naviculo-cuboid bone in the ankle. Tragulids occurred in Africa and Eurasia back to the Miocene, and the more advanced gelocids are known from the upper Eocene and lower Oligocene. At the end of the Oligocene, the first ruminants began to appear with teeth more advanced than those of tragulids. From early in the Miocene they began to be recognizable as giraffes, deer, or antelopes, although the last were relatively uncommon before the late Miocene. Much remains to be learned about the detailed early history of these groups. Several different giraffids lived in later Miocene and early Pliocene times, but the group has since declined to only two species. Deer gradually acquired more complicated antlers, which became very large in some lineages. Different subfamilies of bovids originated in Eurasia and Africa, and it is of zoogeographic interest that representatives of African subfamilies have been found as fossils in northern India and Pakistan.

CLASSIFICATION

Annotated classification. The following classification is principally based on that of American paleontologist George Gaylord Simpson, with alterations in the bovid subfamilies, in the placing of early relatives of giraffes and deer in a giraffoid subfamily Palaeomerycinae, and in the placing of hypertragulids and protoceratids with camels. Groups indicated by the dagger (†) are known only as fossils.

ORDER ARTIODACTYLA

Cloven-hoofed ungulates, the major group of herbivorous mammals. Weight supported mainly through 3rd and 4th toes; astragalus with upper and lower articulations rounded. Stomach compound and, with intestines, enlarged for plant digestion. About 150 species.

Suborder Suiformes

Complete dentition, bunodont (low-cusped) molars, short legs, 4-toed feet are among their important characteristics.

†*Infraorder Palaeodonta*

Eocene to lower Miocene. Primitive, small-brained artiodactyls. Two superfamilies, Dichobunoidea and Entelodontoidea, with 5 and 3 families, respectively, and collectively about 30 genera.

Infraorder Suina

Lower Oligocene to present. Includes the living pigs, peccaries, and their likely ancestors and extinct relatives.

Family Suidae (pigs). Lower Oligocene to present; Old World. Small to moderate size; shoulder height to about 100 cm (39 in.). Coarse hair. Omnivorous, with sharp-edged tusks. Five Recent and about 22 fossil genera.

Family Tayassuidae (peccaries). Differ from pigs by having 1 fewer incisor and premolar, smaller canines, less advanced cheek teeth; hindleg with a cannon bone, more complicated stomach and more densely haired coat.

Infraorder Ancondonta

†*Family Anoplotheriidae*. Eocene and Oligocene; Europe; uncertain relationships.

†*Family Anthracotheriidae*. Eocene to Pleistocene. Mainly Old World, a few in North American Oligocene. Large, with cheek teeth showing beginnings of selenodonty.

Family Hippopotamidae (hippopotamuses). Middle Pliocene to present. Thought to be derived as late as the Pliocene from anthracotheres. Old World, now restricted to Africa. One large (shoulder height to 170 cm [67 in.]; weight to 3,000 kg [6,600

The earliest artiodactyls

Camels and horses in North America

lb]) and one small species (height to 90 cm [35 in.]). Feed on land but frequently resort to water.

†*Family Camotheriidae*. Mainly Oligocene; small European forms of uncertain relationships.

†*Infraorder Oreodontia*

†*Family Merycoidodontidae* (North American oreodonts). Eocene to early Pliocene. No suppression of upper incisors, an incisiform lower canine, selenodont cheek teeth, short faces and short limbs.

†*Family Agriochœriidae*. Eocene to lower Miocene. Close to the above family but with clawed feet.

Suborder Tylopoda

Some reduction of upper incisors, reduced premolars, selenodont cheek teeth; cannon bones present. Feet became 2-toed early in the geological history of the group.

†*Family Hypertragulidae*. Upper Eocene to lower Miocene; North America. Like Old World Tragulina (see below) but with a canine-like first lower premolar. Fused naviculo-cuboid in the ankle.

†*Family Protoceratidae*. Oligocene to lower Pliocene; North America. Some with horns on the top and at the front of the skull. Later ones with hindleg cannon bone. Generally considered close to hypertragulids, but failed to fuse navicular and cuboid bones.

Family Camelidae (camels and lamoids). Upper Eocene to present; now a relict group, represented in southern South America, in Asia, and in North Africa. Red blood corpuscles oval. Gallbladder absent. The hump of the 2 Old World camels is composed of fibrous connective tissue and fat.

†*Family Niphodontidae*. Eocene and lower Oligocene of Europe. Already 2-toed, despite their antiquity, and tentatively placed with the camels.

Suborder Ruminantia (ruminants)

Upper incisors lacking; lower canine incisor-like; cheek teeth selenodont. Fused magnum-trapezoid bone in the wrist. Two-toed feet evolved within suborder.

Infraorder Tragulina

†*Superfamily Amphimerycoidea*

†*Family Amphimerycidae*. European Eocene and Oligocene of Europe; poorly known. Asian *Archaeomeryx*, usually placed in the Hypertragulidae, but may fit here; it retained upper incisors and had a fused naviculo-cuboid in the hind leg.

†*Family Gelocidae*. Eocene-Oligocene of Europe and Asia.

Superfamily Traguloidea

Family Tragulidae (chevrotains). Miocene to present. Sabre-like upper canines in males; incipiently selenodont molars. Bony carapace often develops above the pelvic girdle in males.

Infraorder Pecora

Mostly with horns or antlers and without upper canines. Hollowed auditory bullae. Four-chambered stomach.

Superfamily Giraffoidea

†*Family Palaeomerycidae*. Upper Oligocene to upper Pliocene; North America, Europe, Asia, Africa. Three subfamilies, the Palaeomerycinae (Old World basal stock for giraffes and deer), Blastomerycinae, and Dromomerycinae (the last two New World).

Family Giraffidae (giraffes and okapi). Miocene to present; Old World, now confined to Africa. Living giraffes long-necked and long-legged; the okapi more compact. Giraffes may be up to 5.5 m (18 ft) in total height. Extinct relatives including the large, short-legged, and grotesquely horned sivatheres. Living species have no gallbladder.

Superfamily Cervoidea

Family Cervidae (deer).

Subfamily Moschinae (musk deer). Pleistocene and present; Asia. One recent species, the musk deer (*Moschus moschiferus*) with sabre-like upper canines, gallbladder present (lacking in most other deer).

Subfamily Muntiacinae (muntjacs). Lower Miocene to present; southern Asia. Includes the living muntjacs (*Muntiacus*), with small, 2-pronged antlers above a long, skin-covered base, and tufted deer (*Elaphodus cephalophus*) with tiny antlers. Eurasian fossil genus *Dicrocerus* had larger, 2-pronged antlers. All have large, curved upper canines.

Subfamily Odocoileinae. Pliocene to present. Characterized by the persistence of the lower ends of the lateral metacarpals, thickly haired skin between the hoofs, and in all except the moose a large interdigital gland in at least the hindfoot. The vomer is fused far back posteriorly with the palate in the Amer-

ican deer and in the reindeer. Includes moose (Eurasian elk), reindeer, roe deer, perhaps the Chinese water deer (which has long canines but no antlers), and the deer of North and South America other than the wapiti (American elk), a cervine.

Subfamily Cervinae. Pliocene to present. Top ends of the lateral metacarpals persist. Smooth skin between the hoofs; interdigital glands lacking. Branching of beam of antler differs from that in New World deer; red deer and wapiti, fallow deer, chital, sika, Père David's deer.

Superfamily Bovoidea

Family Antilocapridae (pronghorn and merycodonts). Miocene to present; North America. Smooth, branched horns consisting of hollow sheaths over bony cores; only sheaths shed annually; in fossil Merycodontinae, 1 or more burrs at base of cores. Teeth high-crowned.

Family Bovidae (cattle, sheep and goats, antelopes). Wild cattle, sheep, and goats were ancestral to domestic livestock not differing in any fundamental characters from antelopes. Great variety in horn shape. Many with high-crowned teeth.

Subfamily Bovinae (cattle and some antelopes). Miocene to present. African tribe Tragelaphini, with keeled, spiral horns and not very advanced teeth, includes eland, kudu, nyala, and bushbuck. Tribes Bosclaphini and Bovini, mainly Eurasian, former including the Indian nilgai, the 4-horned antelope, and some extinct forms, the latter including cattle, bison, and buffalo. Bovini are descended from extinct Boselaphini; large size, up to 180 cm (71 in.) at the shoulder. Teeth specialized for grazing. Animals wallow frequently.

Subfamily Cephalophinae (antelopes). Mostly small, with tiny horns set toward the back of the head. Gallbladder lacking. Generally forest-living, African duikers.

Subfamily Hippotraginae (antelopes). Moderate-sized, stocky, mainly grazing antelopes; shoulder height 60–160 cm (24–63 in.). High-crowned teeth. Tribe Hippotragini includes the roan, sable, oryx, and addax antelopes; and Reduncini the reedbuck, kobs, lechwes, and waterbucks. Oryx and addax inhabit arid areas; others near water with adjacent cover or high grass. All except nearly extinct Arabian oryx are now native only to Africa, but the subfamily formerly occurred in India.

Subfamily Alcelaphinae (antelopes). African wildebeests, hartebeests, topis, and several extinct lineages. Long-faced, plains-living, grazing antelopes. Sometimes included in the Hippotraginae.

Subfamily Antilopinae (antelopes). Tribe Neotragini includes some small African antelopes, and Antilopini include gazelles, springbok, and the Indian blackbuck. Graceful, long-legged antelopes of arid, open country. Subfamily may also include the saiga and the Tibetan chiru (tribe Saigini), hitherto classified with Caprinae.

Subfamily Caprinae (sheep, goats, and relatives). Mainly Eurasian. Moderate-sized, shoulder height often 90–100 cm (35–39 in.). High-crowned teeth. Agile animals; many species in mountains or on steep, rocky slopes. Tribe Caprini comprises sheep and goats; Ovibovini the musk-ox (*Ovibos*), the Chinese takin, and some bizarre extinct forms; Rupicaprini the chamois, serow, goral, and Rocky Mountain goat.

Critical appraisal. The great 18th-century classifier Carolus Linnaeus recognized the camels and ruminants as associated but placed some nonartiodactyls with them. It was the French naturalist Henri de Blainville who, at the beginning of the 19th century, first recognized the complete order of artiodactyls as it is accepted today. Nine discrete groups exist among the living forms: pigs, peccaries, hippopotamuses, camels, chevrotains, deer, giraffes, pronghorn, and bovids; their classification presents no great problems, apart from a few genera. Fossils, however, bring confusion to various schemes.

The relationship of North American Tertiary artiodactyls to those of the Old World is a basic question in the study of their zoogeography, history, and classification. It can be agreed that the camels evolved in North America and are as old as the tragulines, which in the Old World were ancestral to ruminants. Most paleontologists today believe that the hypertragulids, protoceeratids, and oreodonts were related to the camels, but others have linked the first two groups with the tragulines and the oreodonts with the anthracotheres. Other questions affecting the higher levels of artiodactyl classification are the placing of the North American native ruminants and whether the earliest Old World pecorans should be taken as giraffoids or cervoids.

The subdivisions of the Bovidae remain controversial. Modifying a classification proposed by German zoologist Max Schlosser, some authorities have grouped the Bo-

vinac, Cephalophinae, and Hippotraginae as the Boödonia and the Alcelaphinae, Antilopinae, and Caprinae as the Aegodontia, to indicate phyletic lines believed to have arisen early in bovid history. Boödonts and aegodonts have evolved differently in Africa and Eurasia, and there is much to be said for developing a classification that reflects these geographical relationships. Until additional evidence of phylogenetic relationships is available, however, the modified version of Simpson's classification above will remain favoured by most taxonomists. (A.W.G.)

BIBLIOGRAPHY

Mammalia. The literature on the biology of mammals is vast. The following list includes only a sample of available sources, with emphasis on those in the English language. Standard general references on the mammals are: F.E. BEDDARD, *Mammalia* (1902); W.H. FLOWER and R. LYDEKKER, *An Introduction to the Study of Mammals, Living and Extinct* (1891); and J.Z. YOUNG, *The Life of Mammals* (1957).

Less technical accounts include: F. BOURLIERE, *Vie et moeurs des mammifères* (1951; Eng. trans., *The Natural History of Mammals*, 1954); and L.H. MATTHEWS, *The Life of Mammals* (1970). Widely used textbooks in mammalogy are: E.L. COCKRUM, *Introduction to Mammalogy* (1962); and D.E. DAVIS and F.B. GOLLEY, *Principles in Mammalogy* (1963). The only thorough up-to-date treatment in English of the families of living mammals is S. ANDERSON and J.K. JONES, JR., (eds.), *Recent Mammals of the World: A Synopsis of Families* (1967). E.P. WALKER *et al.*, *Mammals of the World*, 3 vol. (1964), is a semitechnical work including illustrations of representatives of most living genera; the third volume is a classified bibliography. Perhaps the most comprehensive reference on the biology of mammals is P.P. GRASSE (ed.), *Traité de zoologie*, vol. 16-17 (1967-68, 1955).

The morphology and classification of early mammals is reviewed by J.A. HOPSON and A.W. CROMPTON, "Origin of Mammals," *Evolutionary Biol.*, 3:15-72 (1969); and J.A. HOPSON, "The Classification of Nontherian Mammals," *J. Mammal.*, 51:1-9 (1970). A standard reference on the classification of mammals is G.G. SIMPSON, "The Principles of Classification and a Classification of the Mammals," *Bull. Am. Mus. Nat. Hist.*, 85:1-350 (1945); the evolution of major mammalian groups is treated by A.S. ROMER in *Vertebrate Paleontology*, 3rd ed. (1966), and *Notes and Comments on Vertebrate Paleontology* (1968). See also JASON A. LILLEGRAVEN, *et al.*, (eds.), *Mesozoic Mammals: The First Two-Thirds of Mammalian History* (1979).

Basic references on the status of endangered species of mammals are: F. HARPER, *Extinct and Vanishing Mammals of the Old World* (1945); and G.M. ALLEN, *Extinct and Vanishing Mammals of the Western Hemisphere with the Marine Species of All the Oceans* (1942).

Summaries of specialized topics in mammalogy include: H.T. ANDERSEN (ed.), *The Biology of Marine Mammals* (1969); S.A. ASDELL, *Patterns of Mammalian Reproduction*, 2nd ed. (1964); L.S. CRANDALL, *The Management of Wild Mammals in Captivity* (1964); R.F. EWER, *Ethology of Mammals* (1968); C.P. LYMAN and A.R. DAWE (eds.), *Mammalian Hibernation* (1960); W.V. MAYER and R.G. VAN GELDER (eds.), *Physiological Mammalogy*, 2 vol. (1963-64); A.G. SEARLE, *Comparative Genetics of Coat Colour in Mammals* (1968); GEORGE G. SIMPSON, *Splendid Isolation* (1980), a study of South American mammals.

Some widely known journals that publish papers dealing exclusively with mammals are: *Journal of Mammalogy*, *Mammalia*, and *Säugetierkundliche Mitteilungen* (all issued quarterly), and *Zeitschrift für Säugetierkunde* (6/yr.). Additional technical literature is published by university and public museums and government agencies. Semitechnical and popular accounts of mammals appear in a wide variety of publications of state game departments, museums, and zoological gardens.

Monotremata. The following are popular works on monotremes: H. BURRELL, *The Platypus* (1927); D. FLEAY, *We Breed the Platypus* (1944), and "Observations on the Breeding of the Platypus in Captivity," *Victorian Nat.*, 61:8-14 (1944); M. GRIFFITHS, *Echidna* (1968). Works of a more technical nature are: M.L. AUGEE, E.H.M. EALEY, and H. SPENCER, "Biotelemetric Studies of Temperature Regulation and Torpor in the Echidna, *Tachyglossus aculeatus*," *J. Mammal.*, 51:561-570 (1970); J.H. CALABY, "The Platypus (*Ornithorhynchus anatinus*) and Its Venomous Characteristics," in W. BUCHFEL, E.E. BUCKLEY, and V. DEULOFEU (eds.), *Venomous Animals and Their Venoms*, 1:15-29 (1968); W.K. GREGORY, "The Monotremes and the Palimpsest Theory," *Bull. Am. Mus. Nat. Hist.*, 88:1-52 (1947); K.W. ROBINSON, "Heat Tolerances of Australian Monotremes and Marsupials," *Aust. J. Biol. Sci.*, 7:348-360 (1954); K. SCHMIDT-NIELSEN, T.J. DAWSON, and E.C. CRAWFORD, JR., "Temperature Regulation in the Echidna (*Tachyglossus*

aculeatus)," *J. Cell. Comp. Physiol.*, 67:63-72 (1966); H.M. VAN DEUSEN and G.G. GEORGE, "Results of the Archbold Expeditions, no. 90, notes on Echidnas (Mammalia, Tachyglossidae) of New Guinea," *Am. Mus. Novit.* no. 2383 (1969).

Marsupialia. Among general works are the following: CHARLES L. BARRETT, *Wild Life of Australia and New Guinea* (1954); CHARLES W. BRAZENOR, *The Mammals of Victoria and the Dental Characteristics of Monotremes and Australian Marsupials* (1950); ALBERT S. LE SOUEF and HARRY BURRELL, *The Wild Animals of Australasia, Embracing the Mammals of New Guinea and the Nearer Pacific Islands* (1926); BASIL J. MARLOW, *Marsupials of Australia* (1962); ELLIS TROUGHTON, *Furred Animals of Australia*, 8th ed. rev. (1966); ERNEST P. WALKER *et al.*, *Mammals of the World*, 3 vol. (1964; 2nd ed., vol. 1, 1968); F. WOOD JONES, *The Mammals of South Australia*, 3 vol. (1923-25).

Journal and magazine articles include: J. PEARSON, "Some Problems of Marsupial Phylogeny," *Rep. Meet. Aust. N.Z. Ass. Advmt. Sci.*, 25:71-102 (1947); H.C. REYNOLDS, "Studies on Reproduction in the Opossum (*Didelphus virginiana virginiana*)," *Univ. Calif. Publ. Zool.*, 52:223-284 (1952); G.G. SIMPSON, "The Affinities of the Borhyaenidae," *Am. Mus. Novit.*, no. 1118 (1941); and "The Beginning of the Age of Mammals in South America," *Bull. Am. Mus. Nat. Hist.*, vol. 91, art. 1 (1948); G.H.H. TATE, "On the Anatomy and Classification of the Dasyuridae (Marsupialia)," *Bull. Am. Mus. Nat. Hist.*, vol. 88, art. 3 (1947); "Studies on the Anatomy and Phylogeny of the Macropodidae (Marsupialia)," vol. 91, art. 2 (1948); and "Studies in the Peramelidae (Marsupialia)," vol. 92, art. 6 (1948); D. FLEAY, "Strange Animals of Australia," *Natn. Geogr. Mag.*, 124:388-411 (1963); J.H. CALABY, "Australia's Threatened Mammals," *Wildlife*, 1:15-18 (1963); H.H. FINLAYSON, "Mitchell's Wombat in South Australia," *Trans. R. Soc. S. Aust.*, 85:207-215 (1961); A.G. LYNE, "Australian Mammals," *Aust. Mus. Mag.*, 12:121-125 (1956); J. MCNALLY, "Koala Management in Victoria," *Wildl. Circ. Vict.*, no. 4 (1957); W.E. POOLE and P.E. PILTON, "Reproduction in the Grey Kangaroo, *Macropus kangaroo*, in Captivity," *C.S.I.R.O. Wildl. Res.* 9:218-234 (1964); G.B. SHARMAN, "Studies on Marsupial Reproduction. III. Normal and Delayed Pregnancy in *Setonix brachyurus*," *Aust. J. Zool.*, 3:56-70 (1955); G.B. SHARMAN and J.H. CALABY, "Reproductive Behaviour in the Red Kangaroo, *Megaleia rufa*, in captivity," *C.S.I.R.O. Wildl. Res.*, 9:58-85 (1964); E.M.O. LAURIE and J.E. HILL, "List of Land Mammals of New Guinea, Celebes and Adjacent Islands, 1758-1952," *Br. Mus. Nat. Hist.* (1954).

Insectivora. A. CABRERA, *Genera mammalium*, vol. 1, *Insectivora*, vol. 2, *Galeopithecina* (1919-25), a Spanish classic—a still unsurpassed illustrated review of living insectivora; PETER CROWCROFT, *The Life of the Shrew* (1957), a review of the biology of the common Eurasian *Sorex araneus*; J.F. EISENBERG and E. GOULD, "The Behaviour of *Solenodon paradoxus* in Captivity with Comments on the Behavior of Other Insectivora," *Zoologica*, 51:49-58 (1966); "The Tenrecs: A Study in Mammalian Behavior and Evolution," *Smithson. Contr. Zool.*, no. 27 (1970), a technical but easily read review of the biology of this family; F.G. EVANS, "The Osteology and Relationships of the Elephant Shrews (Macroscelididae)," *Bull. Am. Mus. Nat. Hist.*, 80:85-125 (1942), a technical comparison of elephant shrews, tree shrews, erinacids, and lemuroids; GILLIAN GODFREY and PETER CROWCROFT, *The Life of the Mole (Talpa europaea Linnaeus)* (1960), a readable, semi-popular account of the common European mole; E. GOULD, N.C. NEGUS, and A. NOVICK, "Evidence for Echolocation in Shrews," *J. Exp. Zool.*, 156:19-38 (1964); K. HERTER, "Das Verhalten der Insektivoren," *Handb. Zool.*, vol. 8, sect. 9, pp. 1-50 (1957), a German language review of the behaviour of all insectivores; M.W. LYON, "Treeshrews: An Account of the Mammalian Family Tupaiidae," *Proc. U.S. Natn. Mus.*, vol. 45 (1913); S.B. MCDOWELL, "The Greater Antillean Insectivores," *Bull. Am. Mus. Nat. Hist.*, 115:117-214 (1958), a technical review of the relationships of solenodons to other insectivores.

Chiroptera. G.M. ALLEN, *Bats* (1939, reprinted 1962), a comprehensive view by a naturalist; K.C. ANDERSEN, *Catalogue of the Chiroptera in the Collection of the British Museum*, 2nd ed., vol. 1, *Megachiroptera* (1912), definitive taxonomy of the Megachiroptera; R.W. BARBOUR and W.H. DAVIS, *Bats of America* (1970), comprehensive coverage of North American bats, distribution, natural history, photographs, and taxonomy; A. BROSSET, *La Biologie des Chiroptères* (1966), the most complete, up-to-date review of all bats; R.G. BUSNEL (ed.), *Animal Sonar Systems*, 2 vol. (1967), the most recent important symposium on echolocation and related systems; R.B. DAVIS, C.F. HERREID II, and H.L. SHORT, *Mexican Free-Tailed Bats in Texas* (1962), the most complete published survey of a given species of bats; M. EISENTRAU, *Aus dem Leben der Fledermäuse und Flughunde* (1957), a review of European bats by an important student of bat natural history and physiology; K. FAFRIG and L. VAN

DER PIJL, *The Principles of Pollination Ecology* (1966), the role of bats in pollination; D.R. GRIFFIN, *Listening in the Dark: The Acoustic Orientation of Bats and Men* (1958), the most recent authoritative summary of echolocation; *Echoes of Bats and Men* (1959), a paperback introduction to echolocation; A. NOVICK and N. LFEN, *The World of Bats* (1970), a recent review of bat behaviour and ecology with action photographs; L. VAN DER PIJL, *Principles of Dispersal in Higher Plants* (1969), the role of bats in dispersal of seeds; D.R. ROSEVEAR, *The Bats of West Africa* (1965), an important compilation on the bats of a particular range; H. SAINT GIRONS, A. BROSSET, and M.C. SAINT GIRONS, "Contribution à la connaissance du cycle annuel de la Chauve-souris *Rhinolophus ferrum-equinum* (Schreber, 1774)," *Mammalia*, 33:357-470 (1969), an extensive study of the biology of one European bat species; J. VERSCHUREN, *Ecologie, biologie, et systématique des Chiroptères* (1957), the primary work on the ecology of bats; B. VILLA-R., *Los murciélagos de México* (1966), an important compilation of knowledge on the bats of Mexico; W. WIMSATT (ed.), *The Biology of Bats* (1971), the definitive current work on bat biology.

Primates. *General works:* J.R. and P.H. NAPIER, *A Handbook of Living Primates: Morphology, Ecology and Behaviour of Nonhuman Primates* (1967), a comprehensive, technical book with one or more photographs of each genus; W.C. OSMAN HILL, *Primates: Comparative Anatomy and Taxonomy*, 8 vol. (1953-70), also technical and heavily illustrated, treating the primates in greater depth; W.E. LE GROS CLARK, *The Antecedents of Man*, 3rd ed. (1971), a readable scientific account of nonhuman primates; F. WOOD JONES, *Arboreal Man* (1916), a classic work that contains much useful information; I.T. SANDERSON, *The Monkey Kingdom* (1957), an extensively illustrated popular book; I. DEVORE (ed.), *Primate Behavior* (1965); and C.R. CARPENTER, *Naturalistic Behavior of Nonhuman Primates* (1964), two technical works on primate ethology; W.W. HOWELLS, *Mankind in the Making*, rev. ed. (1967), an informative book on human evolution for the general reader; S. ZUCKERMAN, *Functional Affinities of Man, Monkeys, and Apes* (1933), a technical dissertation on primate structure; IAN TATTERSALL, *The Primates of Madagascar* (1982), a survey of many lower primates; and G.A. DOYLE, *The Study of Prosimian Behavior* (1979).

Tree shrews: W.E. LE GROS CLARK, "The Myology of the Tree-Shrew (*Tupaia minor*)," 1924:461-497; "On the Brain of the Tree-Shrew (*Tupaia minor*)," 1924:1053-1074; "On the Skull of *Tupaia*," 1925:559-567; and "On the Anatomy of the Pen-Tailed Tree Shrew (*Ptilocercus lowii*)," 1926: 1179-1309, all in the *Proc. Zool. Soc. Lond.* (1924-26), four important papers on tree shrews; R.D. MARTIN, "Treeshrews: Unique Reproductive Mechanism of Systematic Importance," *Science*, 152:1402-1404 (1966); and M.W. SORENSON and C.H. CONAWAY, "Observations on the Social Behaviour of Tree Shrews in Captivity," *Folia Primatologica*, 4:124-145 (1966), two technical papers.

Tarsioids: W.E. LE GROS CLARK, "Notes on the Living Tarsier (*Tarsius spectrum*)," *Proc. Zool. Soc. Lond.*, 1924:217-223 (1924); and H.H. WOOLLARD, "The Anatomy of *Tarsius spectrum*," *ibid.*, 1925:1071-1184 (1925), technical papers concerned with the relationships of tarsiers; H. SPRANKEL, "Untersuchungen an *Tarsius*. I. Morphologie des Schwanzes nebst ethologischen Bemerkungen," *Folia Primatologica*, 3:153-188 (1965), a summary, in German, of information on the taxonomy of tarsioids.

Lemuroids: J.J. PETTER, *Recherches sur l'écologie et l'éthologie des Lémuriens malgaches* (1962); "Ecological and Behavioral Studies of Madagascar Lemurs in the Field," *Ann. N.Y. Acad. Sci.*, 102:267-281 (1962); and "The Lemurs of Madagascar," in I. DEVORE (ed.), *Primate Behavior* (1965), technical works on a variety of lemurs; A. JOLLY, *Lemur Behavior* (1967), a readable, semipopular account, based on personal observations; A. WALKER, "Patterns of Extinction Among the Subfossil Madagascan Lemuroids," in P.S. MARTIN and H.E. WRIGHT (eds.), *Pleistocene Extinctions* (1967), a technical discussion of lemur evolution.

New World monkeys: R.I. POCKOCK, "On the External Characters of the South American Monkeys," *Proc. Zool. Soc. Lond.*, 1920:91-113 (1920); and C.R. CARPENTER, "Behavior of Red Spider Monkeys in Panama," *J. Mammal.*, 16:171-180 (1935), two technical articles; L.A. ROSENBLUM and R.W. COOPER (eds.), *The Squirrel Monkey* (1968), a detailed account of the biology of one species; A. CABRERA, *Catalogo de los mamíferos de América del Sur* (1957), a broad treatment, in Spanish, of the taxonomy and distribution of South American monkeys; P. HERSHKOVITZ, "Mammals of Northern Colombia; Preliminary Report No. 4: Monkeys (Primates), with Taxonomic Revisions of Some Forms," *Proc. U.S. Natl. Mus.*, 98:323-427 (1951), a scientific treatment of some South American forms.

Old World monkeys: R.I. POCKOCK, "The External Characters of the Catarrhine Monkeys and Apes," *Proc. Zool. Soc., Lond.*, 1925:1479-1579 (1925); and "The Monkeys of the Genera

Pithecus (or *Presbytus*) and *Pygathrix* Found to the East of the Bay of Bengal," *ibid.*, 1934:895-961 (1934), two articles providing much technical information on the langur; C. G. HARTMAN and W.L. STRAUSS (eds.), *The Anatomy of the Rhesus Monkey, *Macaca mulatta** (1933 and 1961), a detailed account of the morphology of this important species; J.R. and P.H. NAPIER (eds.), *Old World Monkeys* (1970), a technical but comprehensive account of the catarrhines; N.C. TAPPEN, "Problems of Distribution and Adaptation of the African Monkeys," *Curr. Anthropol.*, 1:91-120 (1960), a scientific account of the zoogeography of the catarrhines.

Anthropoid apes: C.F.M. SONNTAG, *Morphology and Evolution of Apes and Man* (1924), a comprehensive and technical book, with many illustrations; R.M. and A.W. YERKES, *The Great Apes* (1929), a classic work, written for the general reader; C.R. CARPENTER, *A Field Study in Siam of the Behavior and Social Relations of the Gibbon (*Hylobates lar*)* (1941, reprinted 1967); G.B. SCHALLER, *The Mountain Gorilla* (1963); and J. VAN LAWICK-GOODALL, *In the Shadow of Man* (1971), three readable accounts of the lives of individual species, based on extensive field work; V. REYNOLDS, *The Apes* (1967); and R. and D. MORRIS, *Men and Apes* (1966), well-illustrated popular books; C.P. GROVES, "Population Systematics of the Gorilla," *J. Zool.*, 161:287-300 (1970); and *Gorillas* (1970), respectively, a technical and popular work on this interesting ape; W.K. GREGORY, *The Anatomy of the Gorilla: The Studies of Henry Cushman Raven, and Contributions by William B. Atkinson and Others* (1950), a carefully illustrated, technical book; TED CRAIL, *Openalk and Whalespeak: The Quest for Interspecies Communication* (1982).

Edentata. The following works, although broad in scope, are particularly rich in information on this group: J.C. BARLOW, "Edentates and Pholidotes," in S. ANDERSON and J.K. JONES, JR. (eds.), *Recent Mammals of the World* (1967), a discussion of the taxonomy of this order and its families; J. DORST, *South America and Central America: A Natural History* (1967), many fine photographs of edentates; E.R. HALL and K.R. KELSON, *The Mammals of North America*, 2 vol. (1959), descriptions and distributional information on North and Central American members of this order; R. HOFFSTETER, "Xenarthra," and R. SABAN, "Palaeonodonta," in J. PIVETEAU (ed.), *L'Origine des Mammifères evolution*, 2 vol., pt. 6 of *Traité de paléontologie* (1958), detailed accounts of early edentates; and E.P. WALKER, *Mammals of the World*, 2nd ed., vol. 1 (1968), an account of the natural history and distribution of every genus in the order.

Lagomorpha. J.N. LAYNE, "Lagomorphs," in *Recent Mammals of the World*, ed. by S. ANDERSON and J.K. JONES (1967), general review of characters, distribution, and habits of living lagomorphs; E.P. WALKER *et al.*, *Mammals of the World*, 2nd ed. (1968), photographically illustrated articles on living genera of lagomorphs; R.M. LOCKLEY, *The Private Life of the Rabbit: An Account of the Life History and Social Behaviour of the Wild Rabbit* (1964), semipopular work dealing with the rabbit in Britain and the disease myxomatosis; H.E. BROADBOOKS, "Ecology and Distribution of the Pikas of Washington and Alaska," *Am. Midl. Nat.*, 73:299-335 (1965), discussion of habits and habitats of North American pikas; M.R. DAWSON, "Lagomorph History and the Stratigraphic Record," in *Essays in Paleontology and Stratigraphy*, ed. by C. TEICHERT and E.L. YOCHELSON (1967), review of fossil lagomorphs and their usefulness for stratigraphic studies; A.A. GUREEV, "Lagomorpha," in *Mammals*, vol. 3, no. 10, *Fauna S.S.S.R.*, Zoologicheskii Institut Akademii Nauk S.S.S.R., New series no. 87 (1964), comprehensive treatment of the morphology and relationships of fossil and living lagomorphs (in Russian); E.R. HALL, "A Synopsis of the North American Lagomorpha," *Univ. Kans. Publ. Mus. Nat. Hist.*, 5:121-202 (1951), on the taxonomy and distribution of North American forms; A.S. LOUKASHKIN, "On the Pikas of North Manchuria," *J. Mammal.*, 21:402-405 (1940), on the behaviour and habitats of Asian pikas.

Rodentia. J.R. ELLERMAN, *The Families and Genera of Living Rodents*, 2 vol. (1940, reprinted 1966), a scientific catalog of existing rodents; C.S. ELTON, *Voles, Mice and Lemmings* (1942, reprinted 1965), a popular, historical record of plagues of various rodents, followed by an analysis of the lemming population fluctuations in northern Labrador; P.P. GRASSE and P.L. DEKEYSER, "Ordre des Rongeurs," in *Traité de zoologie*, vol. 17, pt. 2, pp. 1321-1525 (1955), a section from a zoological reference work dealing with the anatomy, ecology, habits, and classification of rodents, with brief descriptions down to the generic level (in French); T.G. HULL, *Diseases Transmitted from Animals to Man*, 5th ed. (1963), descriptions of diseases, with discussion of causative factors and carriers, many of which are rodents; D. MACCLINTOCK, *Squirrels of North America* (1970), an account of the ecology, habits, and relationships of North American squirrels; H.G.Q. ROWETT, *The Rat as a Small*

Mammal, 2nd ed. (1965), a student laboratory manual with detailed notes on the dissection of the rat; M. SHORTEN, *Squirrels* (1954), a study of the red and gray squirrels in Great Britain for the general reader; G.D. SNELL (ed.), *Biology of the Laboratory Mouse*, 2nd ed. (1966), a standard reference work on all aspects of mouse raising; E.P. WALKER *et al.*, *Mammals of the World*, 3 vol. (1964), descriptions at the generic level of all living mammals, including every recognized genus of rodent, in most cases with pictures; L. WILSSON, *Bäver* (1964; Eng. trans., *My Beaver Colony*, 1968), a popular account of beaver behaviour; and S.E. WOODS, *The Squirrels of Canada* (1980), a guide to 22 species.

Carnivora. For general information on the biology of the carnivores, see E.P. WALKER *et al.*, *Mammals of the World*, 3 vol. (1964), in which each genus is described and illustrated, along with a brief summation of its biology. The taxonomy of carnivores is discussed in G.G. SIMPSON, "Principles of Classification and a Classification of Mammals," *Bull. Am. Mus. Nat. Hist.*, vol. 85 (1945), a classic work on classification, followed by most recent mammalogists; and H.J. STAINS, "Carnivores and Pinnipeds," in S. ANDERSON and J.K. JONES, JR. (eds.), *Recent Mammals of the World: A Synopsis of Families* (1967). F.E. BEDDARD, *Mammalia* (1902), is a definitive early work on mammalian anatomy. The paleontology of the Carnivora is summarized in two works by A.S. ROMER: *Vertebrate Paleontology*, 3rd ed. (1966), fundamental to an understanding of fossil forms, and *Notes and Comments on Vertebrate Paleontology* (1968), containing additional information not found in the general text. Books devoted to particular subgroups of the carnivores include A. DENIS, *Cats of the World* (1964), an excellent summary of the status of all members of the Felidae; C.J. HARRIS, *Otters* (1968), a fine summary of the status of the otters around the world; H.E. HINTON and A.M.S. DUNN, *Mongoose: Their Natural History and Behavior* (1967), which contains much interesting information that is difficult to find elsewhere, except in scattered literature; R.J. HARRISON *et al.* (eds.), *The Behavior and Physiology of Pinnipeds* (1968), an excellent summation of recent knowledge on these aquatic carnivores; and V.B. SCHEFFER, *Seals, Sea Lions, and Walruses: A Review of the Pinnipedia* (1958), a major work on the taxonomy of the Pinnipedia.

Cetacea. G.M. ALLEN, "The Whalebone Whales of New England," *Mem. Bost. Soc. Nat. Hist.*, 8:107-322 (1916), containing much valuable information on the baleen whales; HARALD T. ANDERSEN (ed.), *Biology of Marine Mammals* (1969), a modern review of physiology, anatomy, and behaviour of cetaceans; F.C. FRASER, *Report on Cetacea Stranded on the British Coasts from 1927 to 1932* (1934), ... 1933 to 1937 (1946), and ... 1938 to 1947 (1953), information on North Atlantic cetaceans; *Handbook of R.H. Burne's Cetacean Dissections* (1952), a good source on whale anatomy; RICHARD J. HARRISON and JUDITH E. KING, *Marine Mammals* (1965), a simple guidebook on most aspects of cetaceans; R. KELLOGG, "The History of Whales—Their Adaptation to Life in the Water," *Q. Rev. Biol.*, 3: 29-76, 174-208 (1928); and "A Review of the Archaeoceti," *Publs. Carnegie Instn.* 482 (1936), important papers on the evolution and adaptation of cetaceans; W.N. KELLOGG, *Porpoises and Sonar* (1961), a preliminary synthesis of studies on porpoise echolocation; N.A. MACKINTOSH and J.F.G. WHEELER, "Southern Blue and Fin Whales," *Discovery Rep.*, 1: 257-540 (1929); L. HARRISON MATTHEWS, "The Humpback Whale, *Megaptera nodosa*," *ibid.*, 17:7-92 (1937); and "The Sei Whale, *Balaenoptera borealis*," *ibid.*, 17:183-289 (1938), three papers that deal with economically important species; J.R. NORMAN and F.C. FRASER, *Giant Fishes, Whales, and Dolphins*, new ed. (1948), an excellent review of the cetacea; K.S. NORRIS (ed.), *Whales, Dolphins, and Porpoises* (1966), a readable account of the overall biology of the group; DALE W. RICE and VICTOR B. SCHEFFER, "A List of the Marine Mammals of the World," *Spec. Scient. Rep. U.S. Fish Wildl. Serv. (Fisheries No. 579)* (1968), a systematic listing of all modern cetaceans; P.F. SCHOLANDER, "Experimental Investigations on the Respiratory Function in Diving Mammals and Birds, Hvalråd Skr., 22:1-131 (1940); and P.F. SCHOLANDER and W.E. SCHEVILL, "Counter-Current Vascular Heat Exchange in the Fins of Whales," *J. Appl. Physiol.*, 8:279-282 (1955), technical papers that deal with special adaptations of cetaceans; E.J. SLIJPER, *Die Cetaceen* (1936, in Dutch) and *Walvisen* (1958; Eng. trans., *Whales, 1962*), comprehensive popular books on whales; V.B. SCHEFFER, *The Year of the Whale* (1969), a hypothetical narrative of the life history of the sperm whale, with factual information on other species; A.G. TOMILIN, *Cetacea*, vol. 9 of *Mammals of the U.S.S.R. and Adjacent Countries* (1967; orig. pub. in Russian, 1957); FREDERICK W. TRUE, "Contributions to the Natural History of the Cetaceans: A Review of the Family Delphinidae," *Bull. U.S. Natn. Mus.* 36 (1889), an early synthesis of cetacea that still has much validity; and ERNEST P. WALKER *et al.*, *Mammals of the World*, vol. 2, *Cetacea*, 2nd ed. (1968), photographs and natural history data of most modern cetaceans. LYALL

WATSON, *Sea Guide to Whales of the World* (1981), is a useful reference and field guide to whales, dolphins, and porpoises.

Proboscidea. L.S. DE CAMP, *Elephant* (1964), suitable for both the specialist and nonspecialist alike; RICHARD CARRINGTON, *Elephants: A Short Account of Their Natural History, Evolution and Influence on Mankind* (1958), an excellent study of living as well as extinct forms; P.E.P. DERANIYAGALA, *Elephas maximus, the Elephant of Asia* (1951), with a considerable amount of information on the Asian elephant; and *Some Extinct Elephants, Their Relatives, and the Two Living Species* (1955), reviewed by G. Gaylord Simpson, a remarkable miscellany of elephant lore and observation, ancient and recent; HENRY F. OSBORN, *Proboscidea* (1939), a comprehensive pioneering work; A.S. ROMER, *Vertebrate Paleontology*, 3rd ed. (1966), a *sine qua non* for students of proboscidean evolution; and IVAN T. SANDERSON, *The Dynasty of Abu: A History and Natural History of the Elephants and Their Relatives, Past and Present* (1962), a good all-around work.

Sirenia. G.M. ALLEN, "Extinct and Vanishing Mammals of the Western Hemisphere, with the Marine Species of All the Oceans," *Spec. Publs. Am. Comm. Int. Wildl. Prot.*, no. 11 (1942), provides information on the causes of extinction and the basis for conservation of sirenians; C. BERTRAM, *In Search of Mermaids: The Manatees of Guiana* (1963), is a readable popular account of the biology of manatees; and C.K. and G.C.I. BERTRAM, "The Sirenia as Aquatic Meat-Producing Herbivores," *Symp. Zool. Soc. (London)*, no. 21, pp. 385-391 (1968), contains an extensive bibliography of the literature on the Sirenia.

Perissodactyla. Information on the ungulates as a group is contained in J. SIDNEY, "The Past and Present Distribution of Some African Ungulates," *Trans. Zool. Soc. Lond.*, vol. 30 (1965), a comprehensive account; and J.R. ELLERMAN and T.C.S. MORRISON-SCOTT, *Checklist of Palaeartic and Indian Mammals, 1758 to 1946*, 2nd ed. (1966).

The following papers are concerned with the biology of certain perissodactyls: W. VON RICHTER, "Untersuchungen über angeborene Verhaltensweisen des Schabrackentapirs (*Tapirus indicus*) und des Flachlandtapirs (*Tapirus terrestris*)," *Zool. Beitr. Neue Folge*, 12:67-159 (1966), one of the few detailed studies of tapirs; C.A.W. GUGGISBERG, *S.O.S. Rhino* (1966), a readable general account of the living rhinoceroses, their biology and conservation; H. KLINGEL, "Soziale Organisation und Verhalten freilebender Steppenzebras," and "Soziale Organisation und Verhaltensweisen von Hartman und Bergzebras," *Z. Tierpsychol.*, 24:580-624 and 25:76-88 (1967-68), descriptions of the behaviour and social organization of zebras; and G.G. SIMPSON, *Horses: The Story of the Horse Family in the Modern World and Through 60 Million Years of History* (1961), an interesting account of the natural history and evolution of the Equidae. MYRON J. SMITH, *Equestrian Studies* (1981), is a comprehensive, classified bibliography of more than 4,600 equestrian studies in English.

Artiodactyla. I.W. CORNWALL, *Bones for the Archaeologist* (1956), on the morphology and identification of bones including artiodactyls; J. DORST and P. DANDELDT, *A Field Guide to the Larger Mammals of Africa* (1970), giving summarized descriptions, habits, ecology, and distribution maps for all African artiodactyls; R.F. EWER, *Ethology of Mammals* (1968), a comprehensive text with much information on artiodactyl behaviour; V. GEIST, "The Evolution of Horn-Like Organs," *Behaviour*, 27:175-214 (1966), on the functional implications of horn shape; G.G. SIMPSON, "The Principles of Classification and a Classification of Mammals," *Bull. Am. Mus. Nat. Hist.*, vol. 85 (1945), which forms the basis for most modern classifications of artiodactyls; T. HALTENORTH in W. KUKENTHAL and T. KRUMBACH, *Handbuch der Zoologie*, vol. 8, pp. 1-167, *Klassifikation der Säugetiere: Artiodactyla* (1963), a weighty classification of artiodactyls (in German); V.G. HEPTNER, A.A. NASIMOVIC, and A.G. BANNIKOV (eds.), *Die Säugetiere der Sowjetunion*, vol. 1, *Paarhufer und Unpaarhufer* (1966), a massive work on Eurasian ungulates, most of it dealing with artiodactyls, originally published in Russian in 1961; A. KEAST, "Comparisons of the Contemporary Mammalian Faunas of the Southern Continents," *Q. Rev. Biol.*, 44:121-167 (1969), a review of zoogeography and adaptations, with further references; P.S. MARTIN and H.E. WRIGHT (eds.), *Pleistocene Extinctions: The Search for a Cause* (1967), a collection of essays on Pleistocene extinctions involving artiodactyls; D. MORRIS, *The Mammals* (1965), a general account with illustrations; G.B. SCHALLER, *The Deer and the Tiger* (1967), a study of the life of Indian artiodactyls; C.A. SPINAGE, *The Book of the Giraffe* (1968), much information well presented for the general reader; W.P. TAYLOR (ed.), *The Deer of North America* (1956), on all aspects of the life of North American deer; J.Z. YOUNG, *The Life of Vertebrates*, 2nd ed. (1962), containing a chapter on artiodactyls; and F.E. ZEUNER, *A History of Domesticated Animals* (1963), a useful source for much information difficult to find elsewhere.

Manchester

Manchester remains an important regional city of Great Britain, dominating much of northwestern England, but it has lost that extraordinary vitality and unique influence that made it such a phenomenon of the Industrial Revolution. Manchester was an urban prototype: in many respects it could claim to be the first of the new generation of huge industrial cities created in the Western world during the past 250 years. In 1717 it was merely a market town of 10,000 people, but by 1851 its textile (chiefly cotton) industries had so prospered that it had become a manufacturing and commercial city of more than 300,000 inhabitants, already spilling out its suburbs and absorbing its industrial satellites. By the beginning of the 20th century, salients of urban growth linked Manchester to the ring of cotton-manufacturing towns—Bolton, Rochdale, and Oldham, for example—that almost surround the city, and a new form of urban development, a conurbation, or metropolitan area, was evolving. By 1911 it had a population of 2,350,000. In the following years, however, the pace of growth slowed dramatically. If the 19th century was Manchester's golden age, when it was indisputably Britain's second city, the 20th century was marked by increasing industrial problems associated with the decline of the textile trades (the result of foreign competition and technological obsolescence).

This article is divided into the following sections:

Physical and human geography	460
The landscape	460
The city site	
Climate	
Architecture and the face of the city	
The people	461
The economy	461
Industry	
Trade and transportation	
Administration and social conditions	462
Government	
Education and social services	
Cultural life	462
History	463
Early settlement and medieval growth	463
Evolution of the modern city	463
Bibliography	463

Physical and human geography

THE LANDSCAPE

The city site. Manchester occupies a featureless plain made up of river gravels and the glacially transported debris known as drift. It lies at an elevation of 133 feet (40 metres) above sea level, enclosed by the slopes of the Pennine range on the east and the upland spur of Rossendale on the north. Much of the plain is underlain by coal measures: mining was once widespread but had ceased by the end of the 20th century. Within this physical unit, known as the Manchester embayment, the city's metropolitan area evolved. Manchester, the central city, is situated on the east bank of the River Irwell and has an elongated north-south extent, the result of late 19th- and early 20th-century territorial expansion. In 1930 the city extended its boundaries far to the south beyond the River Mersey, to annex 9 square miles (23 square kilometres) of northern Cheshire. Two large metropolitan boroughs adjoin the city of Manchester on the west and southwest: Salford and Trafford. Together these three administrative units form the chief concentration of commercial employment. From this core, suburbs have spread far to the west and south, chiefly within the historic county of Cheshire. To the north and east of Manchester, smaller industrial towns and vil-

lages, mixed with suburban development, merge into one another and extend as a continuous urban area to the foot of the encircling upland. Close to the upland margin lies a ring of large towns, which were traditionally the major centres of the cotton-spinning industry—Bolton, Bury, and Rochdale to the north and Oldham, Ashton-under-Lyne, and Stockport to the east.

The urban structure of metropolitan Manchester is determined largely by its industrial zones. By far the most important of these is the one bisecting it from east to west. This contains most of the heavier industry—petrochemicals on the Ship Canal near Irlam, electrical engineering in Trafford Park and Salford, and machine tools and metal fabrication in eastern Manchester. Industry in the south is confined to a few compact, largely planned factory estates, notably at Altrincham and Wythenshawe. North and east of Manchester, ribbons of long-established industry follow every railway, river valley, and abandoned canal. The electrochemical industries of the Irwell valley, the dyestuffs of the Irk, and, everywhere, the old textile mills (many converted to new industrial uses) are the dominant features.

Climate. Manchester's climate is most kindly described as mild, moist, and misty. The temperate climate is without extremes: winters are mild, with a January mean temperature of 39° F (4° C), and summers are cool, with a July mean temperature of 59° F (15° C). Occasional high-pressure systems produce cold, clear spells in winter or hot droughts in summer, but these rarely persist. Winds from the west and south prevail, and these bathe the city in frequent gentle rain derived from the almost constant succession of Atlantic weather systems. The annual rainfall, 32 inches (818 millimetres), is not notably high by the standards of western Britain, but it occurs on no less than half of the days in an average year. There is little reliable seasonal variation, but the months of March through May offer the best chance of prolonged dry spells.

The wet Atlantic air banked against the Pennine slopes to the east of the city produces extreme cloudiness; on about 70 percent of the days of the year, the afternoon sky is at least half covered by cloud. This limits sunshine, which was further reduced by air pollution during the decades of the city's industrial prosperity. Up to about 1960 the city centre recorded the abnormally low total of only 970 sunshine hours annually. Foul fogs were another problem of the man-made industrial climate. Manchester then had an average of 55 days of serious fog in a typical year, and the death rate from respiratory diseases surged following these fog episodes. But the city's peculiarly sunless and fog-bound winter climate was transformed by effective air-pollution control. Annual hours of bright sunshine have risen to about 1,300, and serious fogs have been reduced to about 20 days each year. This has been a major factor in reducing the incidence of two formerly endemic diseases, bronchitis and tuberculosis, which had given the city an unenviably high death rate.

Architecture and the face of the city. Manchester's extraordinary 19th-century wealth left a permanent record in an architectural variety and virtuosity that makes the city centre an outdoor museum of styles from Greek classical to early tall steel-framed structures. Commercial firms vied to commission the best architects to design offices and warehouses of ornate splendour, and the public buildings were intended to outshine London's. Thus, banks occupied Greek temples or turreted Gothic castles, and warehouses were given the facades of Venetian palaces. The offices of the Ship Canal Company were given a Grecian colonnade perched high above street level, and the Town Hall, designed by Alfred Waterhouse, is regarded as perhaps the ultimate in Victorian Gothic fantasies.

Conserving this priceless architectural heritage has pre-

sented great problems. Many of the buildings are protected landmarks but are unsuited to modern commercial needs, though some imaginative conversions have taken place. The Royal Exchange, once the hub of the textile trade, contains the old trading floor, the largest room in Europe; it now houses a freestanding theatre-in-the-round. The old Central Station, a huge glazed train shed, has been converted into an exhibition centre. A complex of buildings at Castlefield, including the world's oldest railway station, has been developed as a regional museum of science and industry.

A wave of office redevelopment in the 1960s and '70s added many steel-and-glass structures to the Manchester skyline. One of the earliest is Manchester's tallest building, the Co-operative Insurance Society tower, at 400 feet (122 metres).

As new shopping centres began to develop in outlying areas, the level of retail trade in the city centre suffered. This led to the development of a large enclosed shopping precinct, the Arndale Centre, which contains a significant proportion of the total retail activity in the city centre. As it grew, however, older shopping streets suffered by the shift of businesses, so that parts of the city core have a run-down, half-abandoned appearance; but this is part of the process by which the Victorian central business district is reshaping itself to meet modern needs.

THE PEOPLE

Greater Manchester is one of the world's most compact and crowded metropolitan areas. The overcrowded conditions explain the chief demographic trend of recent years, that of population loss by out-migration. Manchester city itself lost almost one-third of its population to migration between 1961 and 1981, one of the highest rates of migrational loss among all British cities. Natural increase is below the national average, for the migration is chiefly of young families of child-bearing age, leaving an older population in the core cities. Thus overall population decline is serious. This trend is also widespread in the other old industrial towns of the conurbation.

Much of this migration is to suburban areas, though there is also an interregional loss of population to more prosperous areas of Britain, and the "dormitory" districts of the fringes (and especially Cheshire to the south) are growing strongly. Thus, the metropolitan area is decentralizing quickly, and its overall population trend is more favourable than those of its major constituent cities. Total metropolitan population has been virtually stable since 1961, with the low rate of natural increase being entirely offset by net out-migration.

Increasingly, families living in decaying substandard housing have been rehoused. Manchester has exported population to overspill estates at Middleton and Hyde, and Salford families have moved to Worsley. All of these are large schemes, involving population transfers of at least 10,000, and all lie within the metropolitan area. There also has been movement to the New Town project at Warrington, a major development point on the Ship Canal, 18 miles (29 kilometres) west of Manchester. Within the city there has been massive redevelopment. The Hulme scheme of the early 1970s involved the rehousing of a population of almost 60,000.

Like many British cities, Manchester experimented in the 1960s with high-rise housing to accommodate families from the slum clearance zones. In the past, row houses had been the traditional housing form in low-income areas of the inner city, and the new high-rise schemes proved to be a social failure—some were demolished within a decade of construction. The emphasis of city planning was shifted from total clearance and replacement of old housing to its conservation and improvement through Housing Action Areas. Thus, old housing is given new life, and the community is kept together: where new dwellings are built, they are the modern equivalent of the traditional row houses.

The out-migration has been partly counterbalanced by immigration from Commonwealth countries, particularly from the West Indies and the Indian subcontinent. Manchester itself has a multiracial immigrant community,

which is chiefly concentrated in the Moss Side area. Some of the textile towns, too, have attracted Commonwealth immigrants, chiefly Indian and Pakistani textile workers. The metropolitan area as a whole has been one of the main magnets to Commonwealth immigrants in Britain.

THE ECONOMY

Industry. There has long been a contrast between the economies of the core city (Manchester itself, together with the industrial areas of Salford and Stretford) and the textile towns that form the northern and eastern margins of the urban cluster. Until the 1960s the latter had narrowly based economies largely dependent on the textile trade, which still provided more than half the employment of women. The former, however, had an economy of greater diversity: manufacturing was varied (including printing and the production of engineering and electrical products, chemicals, and clothing), and a broad range of service activities gave stability to the economy. This old pattern of contrast was breaking down in the late 20th century, as the core city lost factory employment at a rapid rate and became increasingly dependent on services while the peripheral towns acquired greater industrial diversity and thus a securer (and locally expanding) manufacturing base.

The entire metropolitan area of Greater Manchester has undergone major economic changes. The textile industry has been reduced to a mere vestige of the enormous manufacture that once underpinned the economy of the city. It continues to decline, despite diversification from cotton to man-made fibres and resultant close links with the chemical industry. The surviving mills have been reequipped for high productivity, but this, too, has had the effect of reducing labour demand. The clothing industry has declined with the textile industry but has remained a significant employer of women, chiefly in many small workshops in the inner city. Much more serious has been the sharp contraction of more modern industries that until the 1970s had served as replacements for the old industries. The decline in engineering, one of the main sources of jobs for men, is especially serious. Within the chemical industry the main growth has been in the production of fine chemicals and pharmaceuticals, with research laboratories located in parkland at Alderley, on the southern fringe. The paper and printing industry is stable, reflecting Manchester's status as the second centre, after Greater London, of newspaper production in England.

Manchester's economy has been moving from an industrial to a postindustrial nature. Services have become the chief employers, with the "thinking" rather than the manual services undergoing expansion. Some services, such as transport and distribution, are declining, but the professions, finance and banking, administration, and general personal services are growing with explosive force. Most of these growth points require well-qualified workers: the declining demand for manual skills and the shift to mental skills have caused selective unemployment, which is clearly a persistent social problem.

The conversion of Manchester into a service city is not an entirely new trend, since the city has been the regional capital of northwestern England for two centuries. The process, however, has been quickened by the rapid decline of industry in the inner city. Clearance of the slum tracts and their subsequent redevelopment have removed entire urban districts that once housed many hundreds of small firms. Nearly half of the employment once available in manufacturing in the inner areas has disappeared. In these districts a disadvantaged and ethnically mixed community experiences unemployment rates that are at least twice the city average.

Part of this loss of factory work in the inner city has been the result of the movement of firms to the fringe of the urban area, not only to planned industrial estates but also to the cotton mills left empty by the decline of the textile trades. Hundreds of mills have been converted to other uses, thereby providing the cheap factory-floor space necessary to young and struggling firms, so that the textile towns have in some degree replaced the inner city as an industrial nursery in which it is possible for new firms to become established.

Decline of manufacturing

Commonwealth immigration

Trade and transportation. Apart from its massive volume of retail and wholesale trade, Manchester has a number of distinctions as a regional service centre. It houses a branch of the Bank of England and the Northern Stock Exchange, the headquarters of the Co-operative Wholesale Society, and one of the major provincial crown courts. Its airport at Ringway, 10 miles south of the city, is the leading British terminal outside London in the volume of international traffic handled and in the diversity of both its European and its transatlantic services. Ringway is owned by the city and is the country's second airfreight terminal.

From 1894 to 1986 Manchester was a seaport, with a group of docks at the head of the 37-mile Ship Canal. The growth in the size of shipping, together with changes in the pattern of maritime trade, led to a slow decline in the use of the waterway, and by the mid-1980s the upper parts had been closed to traffic. The lower reaches of the canal remained open and busy, serving the needs of bankside industries, especially the huge oil-refining and chemicals complex at Ellesmere Port. New industrial and commercial uses for the derelict terminal docks have been developed.

Public transport in Greater Manchester is coordinated by a Passenger Transport Executive, and it relies heavily on an integrated system of bus routes. The system faces private competition, however, especially from flexible minibus services. A dense network of commuter rail services, largely electrified, is managed in partnership with British Rail.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. Although the metropolitan area of Greater Manchester is a single cohesive socioeconomic unit, its local government has been fragmented for much of its history. The dominant unit is the metropolitan borough of Manchester, which carries the financial burden of supplying central facilities (major museums and libraries and the airport) for the area as a whole. There are nine other metropolitan boroughs, each independent and able to develop its own social, educational, and planning policies.

The Local Government Act of 1972 (in effect from 1974) created a metropolitan county of Greater Manchester, divided into metropolitan boroughs, including the city of Manchester. The county administered a number of general services (e.g., strategic planning, transport, and recreation), while the boroughs handled the main range of services (e.g., education, housing, and most personal and household services). The metropolitan county of Greater Manchester lost its administrative powers in 1986, however. Some of the general services that it had provided were

taken over by specialist successor authorities, but many of its administrative powers passed to the city of Manchester and the other individual metropolitan boroughs, which are in effect now unitary authorities.

Education and social services. Of all Manchester's pioneer cultural achievements, none has prospered more than the Victoria University of Manchester. After its foundation in 1851 at a site in Quay Street, the college received a charter in 1872 and began growth on its present site in 1873. By 1880 it had combined with member colleges in both Leeds and Liverpool to form a federal institution. Since becoming a separate body again in 1903, the university has grown to become one of the largest in Britain. The faculty of technology has become autonomous as an Institute of Science and Technology, and, with the establishment of the University of Salford in 1967 and the growth of a large polytechnic, there are now four institutions of higher learning in and near the city.

The city provides the complete range of social and welfare services within the British system, but its special strength lies in health services and medical education. The Victoria University of Manchester has the largest medical school in western Europe; it is linked to three large groups of teaching hospitals that provide specialist treatment. One of the most distinguished of these is the Christie Hospital, a major centre for cancer research.

CULTURAL LIFE

The cultural life of Manchester suffered some losses during the 20th century. For example, its prestigious newspaper, *The Guardian*, has (in the Mancunian view) fled to London and dropped the city's name from its title. However, the Lowry, an architecturally innovative centre for the visual and performing arts, opened in 2000 and signaled the city's cultural revival at the beginning of the 21st century. Music maintains its strength. The Hallé concerts reached their centenary in 1958, and the orchestra continues to maintain its international reputation.

The city has a large number of private, public, and specialized libraries. The municipal library, with more than 25 branches, has its headquarters at St. Peter's Square. Manchester also houses the notable John Ryland University Library (now part of the Victoria University of Manchester library) and Chetham's Library, one of the first free public libraries in Europe.

Among the galleries and museums, the Whitworth Art Gallery and the Manchester City Art Gallery are particularly well known. The latter contains a fine collection of paintings, sculpture, silver, and pottery and is supplemented by several branch galleries. The Manchester Museum has special exhibits of Egyptian and Japanese objects, as

Greater
Manchester

Libraries



The Guardian, Manchester and London

Manchester skyline viewed from the mathematics building of Manchester University Institute of Science and Technology.

well as natural history collections and an aquarium. The Museum of Science and Industry highlights Manchester's industrial heritage.

There are two major association football (soccer) clubs, including Manchester United, one of the world's most famous and popular teams. At the grounds of the Old Trafford Cricket Club, test matches are played against overseas cricket teams visiting Great Britain. Manchester also has an active pop music scene, which revolved around Factory Records in the 1980s and has given rise to several influential rock bands, including Joy Division, the Smiths, and Oasis.

History

EARLY SETTLEMENT AND MEDIEVAL GROWTH

Early in the Roman conquest of Britain, a fort was established (AD 78–86) on a low sandstone plateau at the confluence of the Rivers Medlock and Irwell. In its first form, the fort was a simple field fortification of shallow ditches, earth banks, and timber palisades. By the early 3rd century, it had been rebuilt in stone and contained a number of buildings; excavations have uncovered evidence of substantial activity. A *vicus* (Latin: "row of houses") of merchants and craftsmen had grown outside the walls, along the well-made road to York. But Roman occupation left no permanent imprint, except to give the modern city its name, derived from Mamucium ("Place of the Breastlike Hill"). There is no evidence of occupation after the 4th century, and the site seems to have lain empty for 500 years. In 919 the West Saxon king Edward the Elder sent a force to repair the Roman site as a defense against the Norsemen, and some traces of this re-occupation have been discovered. By then, however, the growth of Manchester had recommenced almost a mile from the fort, at the junction of the Rivers Irk and Irwell near the present cathedral.

The Norman barony of Manchester was one of the largest landholdings in Lancashire, and its lords built a fortified hall close to the church. During the 13th century, Manchester began its transition from village to town, and sometime before 1301 a charter was granted. Although Manchester was acquiring regional importance, it was subordinate to its near neighbour, Salford, which was the capital manor of the hundred (district) and which had an earlier borough charter. The full development of the medieval borough followed the establishment in 1421 of a college of priests to take charge of the church. Part of the college survives as Chetham's Hospital, while a free church school set up in 1506 became the Manchester Grammar School in 1515, founded by Hugh Oldham, bishop of Exeter.

EVOLUTION OF THE MODERN CITY

By the 16th century Manchester was a flourishing market borough important in the wool trade, exporting cloth to Europe via London. By 1620 a new industrial era had begun with the weaving of fustian, a cloth with a linen warp but a cotton weft. This was the origin of the cotton industry that was to transform southern Lancashire after 1770. As the trade grew, Manchester expanded and "improvements" were added, including the fine square and church of St. Ann (1712).

From the 1760s onward, growth quickened with the onset of the Industrial Revolution. The first canal, bringing cheap coal from Worsley, reached the town in 1762; later extended, it linked Manchester with the Mersey and Liverpool by 1776 and so served the import-export needs of the cotton industry. Manchester's first cotton mill was built in the early 1780s. By 1800 Manchester was said to be "steam mill mad," and by 1830 there were 99 cotton-spinning mills. The world's first modern railway, the Liverpool and Manchester, was opened in 1830, and by the 1850s the greater part of the present railway system of the city was complete. Despite its growth to a population of more than 70,000 by 1801, the town had no system of government and was still managed, like a village, by a manorial court leet (a court held semiannually by the lord of the manor or his steward to conduct local gov-

ernment). A police force was established in 1792, but not until 1838 did a charter of incorporation set up an elected council and a system of local government.

Manchester's economic history during the second half of the 19th century was one of growth and diversification. The city became less important as a cotton-manufacturing centre than as the commercial and financial nucleus of the trade; on the floor of the Royal Exchange, the yarn and cloth of the entire industry was bought and sold. From an early textile-machinery industry, many specialized types of engineering developed. Products included steam engines and locomotives, armaments, machine tools, and, later, those of electrical engineering. The opening of the 37-mile Manchester Ship Canal (1894) linked Manchester, via the Mersey estuary at Eastham, to the Irish Sea and the world markets beyond. By 1910 Manchester had become the fourth port of the country, and alongside the docks, at Trafford Park, the first (and still the largest) industrial estate in Britain was developed. New industries also took sites there, with a prominent role played by such American companies as Westinghouse and Ford, the latter moving to Essex in 1929. At its height, more than 50,000 workers were accommodated within factories of the estate, though that number later declined.

The Manchester of the 19th century was a city of enormous vitality not only in its economic growth but also in its political, cultural, and intellectual life. *The Manchester Guardian* became Britain's leading provincial newspaper, achieving international influence, while the Hallé Orchestra was its equal in the world of music. Owens College (now known as Victoria University of Manchester) became the nucleus of the first and largest of the great English civic universities, while the academic success of the Manchester Grammar School made it something of a model in the development of selective secondary education in England. Politically, Victorian Manchester often led the nation: in the agitation for parliamentary reform and for free trade, its influence was crucial. The Peterloo Massacre of 1819 arose from a peaceful political assembly, held on fields near the city, to demand parliamentary reform. In the period 1842–44 the German social philosopher Friedrich Engels lived in Manchester, and his influential book *Condition of the Working Class in England* (1845) was based on his experiences there. Among its other intellectual achievements were John Dalton's development of the atomic theory as the foundation of modern chemistry and the work of the "Manchester school" in the application of economic principles to the problems of commerce, industry, and government.

There was a price to be paid for this precocious growth. In its urban fabric, inner Manchester remained essentially a 19th-century city, and by the late 20th century it faced massive redevelopment problems. An industrial collar of obsolescent factory zones encircled the city centre, and huge areas of old slum housing survived with little renewal into the 1960s. Although the city centre was devastated by an Irish Republican bomb in 1996, much of it was rebuilt by the beginning of the 21st century. Manchester, then, is a city in transition: its face is being transformed by redevelopment, and its dependence on the insecure base of the textile industries is declining with the growth of a much broader economic structure.

BIBLIOGRAPHY. Historical studies include ALAN J. KIDD and K.W. ROBERTS (eds.), *City, Class, and Culture: Studies of Social Policy and Cultural Production in Victorian Manchester* (1985); JOHN H.G. ARCHER (ed.), *Art and Architecture in Victorian Manchester* (1985); GARY S. MESSINGER, *Manchester in the Victorian Age* (1985); NICHOLAS J. FRANGOPULO (ed.), *Rich Inheritance: A Guide to the History of Manchester* (1962, reissued 1969); D.A. FARNIE, *The Manchester Ship Canal and the Rise of the Port of Manchester, 1894–1975* (1980); and SIENA D. SIMON, *A Century of City Government: Manchester, 1838–1938* (1938). For Manchester in its regional context, see C.F. CARLIER (ed.), *Manchester and Its Region* (1962); T.W. FREEMAN, H.B. RODGERS, and R.H. KINVIG, *Lancashire, Cheshire and the Isle of Man* (1966); L.P. GREEN, *Provincial Metropolis* (1959); and H.P. WHITE (ed.), *The Continuing Conurbation: Change and Development in Greater Manchester* (1980). (H.B.Ro./Ed.)

Early
growth

The
industrial
era

19th-
century
achievement

Manila

Manila, the national capital and chief city of the Republic of the Philippines, is the centre of the country's economic, political, social, and cultural activity. Located on Luzon Island, it spreads along the eastern shore of Manila Bay at the mouth of the Pasig River. The city's name, originally Maynilad, is derived from that of the *nilad* plant, a flowering shrub adapted to marshy conditions, which once grew profusely along the banks of the river; the name was shortened first to Maynila and then to its present form. The city proper encompasses an area of approximately 15 square miles (38 square kilometres). In 1975, by presidential decree, Manila and its contiguous cities and municipalities were integrated to function as a single administrative region, known as Metropolitan Manila (also called the National Capital Region), with an area of 246 square miles.

Manila has been the principal city of the Philippines for four centuries, and it is the centre of its industrial development as well as the international port of entry. It is situated on one of the finest sheltered harbours in the region, about 700 miles (1,100 kilometres) southeast of Hong Kong. The city has undergone rapid economic development since its destruction in World War II and its subsequent rebuilding; it is now plagued with the familiar urban problems of pollution, traffic congestion, and overpopulation. Measures have been taken, however, to ameliorate these problems.

This article is divided into the following sections:

Physical and human geography	464
The landscape	464
The city site	
Climate	
Plant and animal life	
The city layout	
The people	465
The economy	466
Industry	
Commerce and finance	
Transportation	
Administration and social conditions	466
Government	
Public utilities	
Health and security	
Education	
Cultural life	467
History	467
Bibliography	467

Physical and human geography

THE LANDSCAPE

The city site. Manila lies on the eastern shore of Manila Bay, a large inlet with access to the sea through a channel 12 miles wide to the southwest. It occupies the low, narrow deltaic plain of the Pasig River, which flows northward to Manila Bay out of a large lake, Laguna de Bay, southeast of the city. Manila Bay lies to the west, the swampy delta of the southward-flowing Pampanga River to the north, the mountains of the Bataan Peninsula to the west, and Laguna de Bay to the southeast. Although the city's area is constricted, it is an excellent port site because of its sheltered harbour, its access by river to inland agricultural areas, and its proximity to the Asian mainland.

Climate. The city is protected from extreme weather conditions by the hills of the Eastern Cordillera to the east and by the mountains of the Bataan Peninsula, which lies west of Manila Bay. The tropical climate is characterized by a wet season that lasts from June to November and by a dry season lasting from December to May. High humid-

ity and thunderstorms are common in July, August, and September, when more rain is received than in other months. The average annual rainfall totals about 82 inches (2,080 millimetres). There is little monthly variation from the mean annual temperature of 81° F (27° C).

Plant and animal life. The landscape of the city of Manila is dotted with tropical trees, including the palm, banyan, and acacia; bamboo thrives in the public parks. Water buffalo, horses, dogs, pigs, and goats are common. The wealth of birdlife includes shrikes, doves, and pigeons, and Manila Bay abounds with sardines, anchovies, mackerel, tuna, snappers, and barracuda. The city's natural beauty is marred, however, by air and water pollution caused by the expansion of industry and the growing number of motor vehicles.

The city layout. The city is bisected by the Pasig River and divided into four administrative divisions comprising 14 districts. The districts developed from the original fortress city of Intramuros ("Within Walls") and the 13 villages located outside its walls. The districts of Tondo, Santa Mesa, Binondo, Santa Cruz, Quiapo, San Miguel, and Sampaloc lie to the north of the river, and Intramuros, Ermita, Malate, Paco, Pandacan, and Santa Ana are to the south. The two sections of the city are connected by several bridges.

Business areas are widespread, but Makati City is the chief centre of trade and commerce. San Miguel is the site of Malacañang Palace, the presidential residence; several universities are located in Sampaloc. The shantytown district of Tondo on the northern shore is the site of Manila North Harbor, the local port, while the international port, Manila South Harbor, is on the southern shore. Intramuros is renowned for its 16th-century San Agustin church and for the ruins of its old walls and of Fort Santiago. Ermita and Malate are choice residential districts and the sites of hotels and embassies; Paco, Pandacan, and Santa Ana are middle-income residential areas.

Metropolitan Manila includes the cities of Manila and Caloocan City to the north, Quezon City to the northeast, and Pasay City (near the shore of Manila Bay) to the south and 13 municipalities. The municipalities include Makati, Mandaluyong, San Juan, Las Piñas, Malabon, Navotas, Pasig, Pateros, Parañaque, Marikina, Muntinlupa, Tagig, and Valenzuela. Metropolitan Manila was created to provide integrated services such as water supply, police and fire protection, and transport and to permit central planning for unified development.

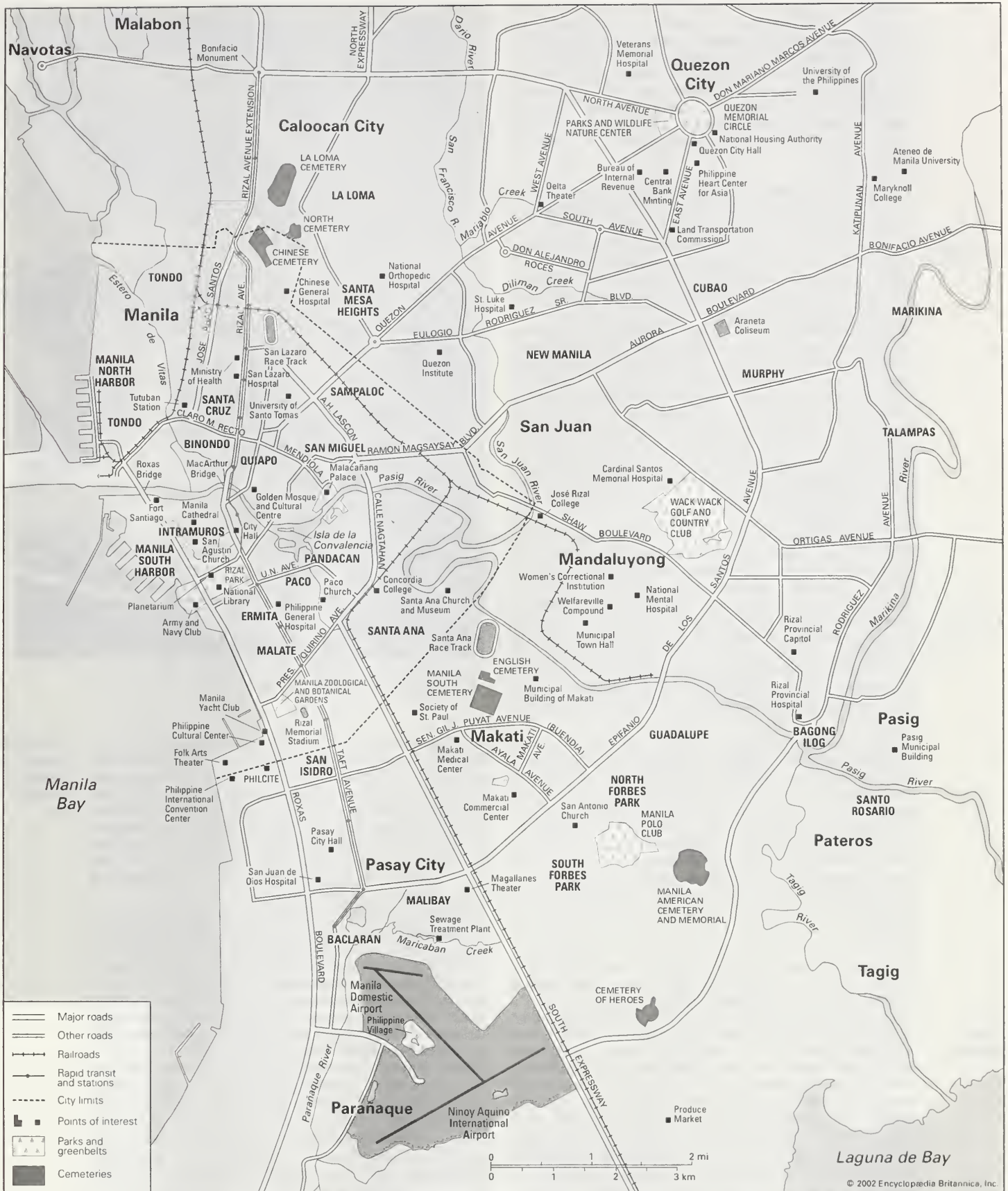
Housing. The city has a chronic housing shortage, and tenement housing projects have been constructed by the government to help house the poor. To provide homes for squatters, the government also has developed resettlement projects around Manila that are easily accessible by land motor transportation.

Residential buildings include the single-family dwelling; the duplex for two independent households; the *accessoria*, whose dwelling units have individual entrances from the outside; the apartment building with common entrance; and the *barong-barong*, a makeshift shack built of salvaged materials (flattened tin cans, scrap lumber, cartons, or billboards) that is common in the poor area of Tondo.

Architecture. Architectural styles reflect American, Spanish, Chinese, and Malay influences. Rizal Park and a number of government buildings were designed by the American architect and city planner Daniel H. Burnham. Modern buildings—including multistoried commercial houses and public and private buildings—are commonly made of reinforced concrete and hollow cement blocks. Houses of modern design—especially low, sprawling ranch houses with spacious lawns—are common in the districts of Ermita and Malate. Spanish-style houses, with tiled roofs, barred windows, and thick walls, were common be-

City districts

Types of dwellings



Manila with its adjoining municipalities.

fore World War II and have remained popular. The churches of the city are American, Spanish, or European in character. The Manila cathedral was rebuilt in the 1950s and is an important landmark. It succeeds five earlier cathedrals—the first dating from the mid-16th century—that were destroyed either by earthquakes or during wartime.

THE PEOPLE

Metropolitan Manila is densely populated and contains a large proportion of the country's population as a result of rural-urban migration. The strain on municipal services has an adverse effect on the quality of life in the urban area. To ameliorate the situation, the government has adopted a "back-to-the-farm" policy and has established

resettlement projects outside of the central urban area.

Almost all the residents of Manila are Filipinos. The largest single foreign community, representing less than 10 percent of the population, is made up of Chinese. The population of the city is predominantly Roman Catholic, although there are some Protestants, Muslims, and Buddhists. The two national Christian churches—the Iglesia ni Kristo and the Philippine Independent, or Aglipayan, Church—have small congregations. (D.C.S.)

THE ECONOMY

Industry. Manila's diverse manufactures include textiles and clothing, publishing and printing, food processing, and the production of tobacco products, paints, drugs, aluminum articles, rope and cordage, shoes, coconut oil, soap, and lumber. The manufacture of electronic products and components has increased. Factories vary in size from small plants located in congested districts to large facilities on the outskirts of the metropolitan area.

Commerce and finance. Metropolitan Manila is the centre of commerce and finance in the Philippines. Makati City is the most important concentration of commercial and financial activity within the metropolitan area. Ayala Avenue is lined with the high-rise headquarters of multinational and Philippine banks, industrial enterprises, and trading houses. Makati City is also one of the most important retail and shopping areas, with its many department stores, enormous shopping malls, specialty stores, and boutiques. Just northeast of Makati City, Mandaluyong, with its new shopping complexes, is emerging as another important retail district.

Trade flourishes within the metropolitan area and between it and the provinces and other countries. Most of the Philippines' imports and exports pass through the port of Manila. Financial institutions headquartered in Metropolitan Manila include such establishments as the Development Bank of the Philippines, the Philippine National Bank, the Philippine Veterans Bank, the Government Service Insurance System, the Social Security System, and many private commercial and developmental banks. Private insurance companies and the Manila Stock Exchange also contribute to the mobilization of savings for investment.

Transportation. Within the region, public transportation is provided principally by buses, jeepneys (small buses built on the chassis of jeeps), and taxis. An elevated rail line, linking Caloocan City and the city of Baclaran (to the south of Pasay City), was completed in 1984. It was the first phase of a transit system, called Light Rail Transit, that eventually will extend throughout the metropolitan area. A new line that will encircle the metropolitan area is under construction. Bus services operate routes to northern and southern Luzon; railroad services operated by the Philippine National Railways also connect the city with northern and southeastern Luzon. Interisland and international transportation is provided by shipping and by domestic and foreign airlines. Manila Domestic Airport and Ninoy Aquino International Airport are both located about 6 miles (10 kilometres) south of the city centre.

Traffic congestion is serious, traffic tending to pile up at the bridges during the morning and evening rush hours. Adjacent towns serve as dormitory suburbs, and many people commute to the city, adding to the traffic problem.

The main international port is Manila South Harbor, enclosed by a low breakwater. Several heavy industries that depend upon imported raw materials are located within the port area. The piers of the local port, Manila North Harbor, are congested with heavy traffic from all ports in the Philippines. It has several warehouses for storage of goods and equipment. Additional port facilities for international shipping have been built, partially on reclaimed land, in the area between the two harbours. (D.C.S./Ed.)

ADMINISTRATION AND SOCIAL CONDITIONS

Government. With the organization of Metropolitan Manila in 1975, the mayor-council type of government was superseded by a manager-commission type. Metropolitan Manila, which is vested with the powers and attributes of a corporation (including the power to make



Monument (centre) to José Rizal, national hero of the Philippines, in Rizal Park, Manila.

Tautomu Nakayama—Shostal Assoc.

contracts, sue and be sued, hold, transfer, and dispose of property, and similar powers), is administered by the Metropolitan Manila Commission. It is headed by a governor as the general manager, who executes policies and measures approved by the commission and is responsible for the discharge of management functions. The other members are the vice governor, as deputy general manager, and three commissioners or board members, one each for planning, finance, and operations. All members are appointed by the president of the Philippines.

Supplanting the city or municipality boards or councils is the Sangguniang Bayan (Municipal Assembly), created for each city or municipality, which helps the metropolitan government in administration and legislation. It is composed of the mayor, vice mayor, councillors, captains of *barangays* (neighbourhoods), and representatives from other sectors, appointed by the president upon recommendation of the local unit. The mayors are the presiding officers in their respective areas of jurisdiction.

Public utilities. Potable water comes from a supply network managed by the Metropolitan Waterworks and Sewerage System. Satisfactory sanitation conditions are maintained by surveillance of markets, restaurants, motion-picture houses, recreation halls, and slaughterhouses. Insecticides are sprayed on open sewers, uncollected garbage, and standing water; garbage is collected by a fleet of trucks that operate night and day. Moreover, workers maintain cleanliness in the metropolitan area and are also responsible for the beautification of the city as directed by the governor of Metropolitan Manila.

Health and security. Health facilities in Manila are among the best in the region. The city government maintains numerous health centres as well as San Lazaro Hospital, where patients are treated free of charge, and subsidizes a number of government hospitals. There are also many missionary and private hospitals in the city.

The Metropolitan Manila Commission

Traffic problems

Police and fire services

Police and fire services are well organized, improved equipment and techniques are used, and personnel are comparatively well paid. To intensify the campaign against crime, more police precincts were created, and the Police Community Relations District was organized. *Barangay* brigades and *barangay tanods* (guards) were also created in each *barangay* of every city and municipality in Metropolitan Manila. They are volunteers and selected leaders of the *barangay* whose duty is to maintain peace and order in their community.

Education. More than 96 percent of the population of 10 years of age or older is literate. More than 100 free public schools are maintained, in addition to the night vocational and secondary schools and the city-supported University of Manila. Educational opportunities are also provided for handicapped children, orphans of school age, and adults. As the education centre of the country, Metropolitan Manila houses many of the major institutions of higher education of the country, including the University of the Philippines (with its main campus in Quezon City), the Philippine Normal College, and the Polytechnic University of the Philippines. There are several universities sponsored by religious bodies, including the University of Santo Tomas (founded in 1611), as well as nonsectarian institutions such as the University of the East and the Far Eastern University.

Manila's universities

CULTURAL LIFE

The centre of the performing arts in the country is the Philippine Cultural Center. There are also the Folk Arts Theater, facing Manila Bay, the renovated historic Metropolitan Theatre, and an open-air theatre in Rizal Park. The many libraries and museums include the National Library and the National Museum, known for its anthropological and archaeological exhibits; the National Institute of Science and Technology, with a scientific reference library and large collections of plants and animals; the geological museum of the Bureau of Mines and Geosciences; the Planetarium; Fort Santiago, which houses original works of the Philippine patriot José Rizal; and the Kamaynilaan (Manila City) Library and Museum, which contains valuable carvings, paintings, and archives.

The city's parks

The foremost outdoor recreational area is Rizal Park, with a Japanese garden, a Chinese garden, an open-air theatre, a playground, a grandstand, and a long promenade adjacent to Manila Bay. Other areas include the Manila Zoological and Botanical Gardens, the Mehan Garden, and Paco Park. Athletic facilities include the Rizal Memorial Stadium and the Jai-Alai Fronton, both located in Manila, and the Araneta Coliseum in Quezon City. Annual festivals and carnivals are held in the sunken garden fronting the City Hall of Manila.

The interweaving of indigenous, Muslim, Spanish, and American influences has given Manila a variegated culture more Western than that of most other Asian countries. The Spanish legacy is represented in the importance of Roman Catholicism in daily life and in the numerous annual fiestas, many of which commemorate patron saints. Folk art in Manila has many expressions, but among its most prominent manifestations are jeepneys, which date to the World War II jeeps left behind by U.S. soldiers. Although jeepneys are a main form of public transportation, they also serve as rolling art galleries that display colourful statues, designs, and flashing lights.

Sports are an important part of daily life in Manila. Basketball, introduced by Americans in the early 1900s, is the most popular team activity. The games of the country's professional teams have high attendance rates, and many enjoy playing basketball in public parks and elsewhere. Golf, tennis, water sports, boxing, cockfights, and Filipino martial arts (arnis, kali, and eskrima) are other activities that have a following among Manilans.

History

In the late 16th century, Manila was a walled Muslim settlement whose ruler levied customs duties on all commerce passing up the Pasig River. Spanish conquistadors under the leadership of Miguel López de Legazpi—first Spanish

governor-general of the Philippines—entered the mouth of the river in 1571. They destroyed the settlement and founded the fortress city of Intramuros in its place. Manila became the capital of the new colony. Outside the city walls stood some scattered villages, each ruled by a local chieftain and each centred on a marketplace. As Spanish colonial rule became established, churches were built near the marketplaces, where the concentration of population was greatest. Manila spread beyond its walls, expanding north, east, and south, linking together the market-church complexes as it did so.

The propagation of Roman Catholicism began with the Augustinian friar Andrés de Urdaneta, who accompanied the expedition of 1571. He was followed by Franciscan, Dominican, Jesuit, and Augustinian priests who founded churches, convents, and schools. In 1574 Manila was baptized under the authorization of Spain and the Vatican as the "Distinguished and Ever Loyal City" and became the centre of Catholicism as well as of the Philippines. At various periods Manila was threatened, and sometimes occupied, by foreign powers. It was invaded by the Chinese in 1574 and raided by the Dutch in the mid-17th century. In 1762, during the Seven Years' War, the city was captured and held by the British, but the Treaty of Paris (1763) restored it to Spain. It was opened to foreign trade in 1832, and commerce was further stimulated by the opening of the Suez Canal in 1869.

The Manila area became the centre of anti-Spanish sentiment in the 1890s, and the execution of Filipino patriot José Rizal in the city in December 1896 sparked a year-long insurrection. During the Spanish-American War the Spanish fleet was defeated at Manila Bay on May 1, 1898, and on August 13 the city surrendered to U.S. forces. It subsequently became the headquarters for the U.S. administration of the Philippines.

The U.S. period was one of general social and economic improvement for the city. U.S. policy encouraged gradual Filipino political autonomy, and to help achieve this goal public schools were established in Manila and throughout the archipelago. The University of the Philippines, founded in 1908, became the apex of the educational system. The city developed into a major trading and tourist centre.

Upon the outbreak of World War II, Manila was declared an open city and was occupied by the Japanese in January 1942. The city suffered little damage during the Japanese invasion but was levelled to the ground during the fight for its recapture by U.S. forces in 1945.

Manila was in shambles when in 1946 it became the capital of the newly independent Republic of the Philippines. The city was rapidly rebuilt, however, with U.S. aid. A significant change in its appearance was brought about by industrialization. In 1948 suburban Quezon City was chosen as the site of a new national capital, but in 1976 Manila again became the capital and the permanent seat of the national government. (D.C.S.)

BIBLIOGRAPHY. General discussions on the history, government, economy, and living conditions in Manila include TEODORO A. AGONCILLO and MILAGROS C. GUERRERO, *History of the Filipino People*, 8th ed. (1990); ALFONSO J. ALUIT, *The Galleon Guide to Manila*, 3rd rev. ed. (1973); ENRIQUE L. VICTORIANO, *Historic Manila: Commemorative Lectures, 1993–1996* (1997); MIGUEL ANSELMO BERNAD, *The Western Community of Manila: A Profile* (1974); ROBERT E. IUKE, *Shadows on the Land: An Economic Geography of the Philippines* (1963); LUNING B. IRA, *Streets of Manila* (1977); FUND FOR ASSISTANCE TO PRIVATE EDUCATION, *The Philippine Atlas*, vol. 1, *A Historical, Economic and Educational Profile of the Philippines* (1975); and DOMINGO C. SALITA, *Geography and Natural Resources of the Philippines* (1974). Topics in military history are discussed in JAMES H. NELSON, *Threshold of Empire and the Battle for Manila: 1898–1899* (1998); R.M. CONNAUGHTON, *The Battle for Manila* (1995); and ALPHONSO J. ALUIT, *By Sword and Fire: The Destruction of Manila in World War II: 3 February–3 March 1945* (1994).

Issues associated with Metropolitan Manila are addressed in MANUEL A. CAOLI, *The Origins of Metropolitan Manila: Social and Political Analysis* (1999); and ERHARD BERNER, *Defending a Place in the City: Localities and the Struggle for Urban Land in Metro Manila* (1997). Specific treatment on population, housing, charter provisions, needs, and resources of the capital city may be found in *Philippine Yearbook* (annual). (D.C.S./Ed.)

The Spanish period

The U.S. period and independence

Mao Zedong

As China emerged from a half century of revolution as the world's most populous nation and launched itself on a path of economic development and social change, Mao Zedong (Wade-Giles: Mao Tse-tung), its principal revolutionary thinker and for many years its unchallenged leader, occupied a critical place in the story of the country's resurgence. To be sure, he did not play a dominant role throughout the whole struggle. In the early years of the Chinese Communist Party, he was a secondary figure, though by no means a negligible one, and even after the 1940s (except perhaps during the Cultural Revolution) the crucial decisions were not his alone. Nevertheless, looking at the whole period from the foundation of the Chinese Communist Party in 1921 to Mao's death in 1976, one can fairly regard Mao Zedong as the principal architect of the new China.

Eastfoto



Mao Zedong, 1966.

Early years. Born on Dec. 26, 1893, in the village of Shao-shan, Hunan Province, Mao was the son of a former poor peasant who had become affluent as a farmer and grain dealer. He grew up in an environment in which education was valued only as training for keeping records and accounts. From the age of eight he attended his native village's primary school, where he acquired a basic knowledge of the Confucian Classics. At 13 he was forced to leave and begin working full time on his family's farm. Rebelling against paternal authority, Mao left his family to study at a higher primary school in a neighbouring county and then at a secondary school in the provincial capital, Ch'ang-sha. There he came in contact with new ideas from the West, as formulated by such political and cultural reformers as Liang Ch'i-ch'ao and the Nationalist revolutionary Sun Yat-sen. Scarcely had he begun studying revolutionary ideas when a real revolution took place before his very eyes. On Oct. 10, 1911, fighting against the Manchu dynasty broke out in Wu-ch'ang, and within two weeks the revolt had spread to Ch'ang-sha.

Enlisting in a unit of the revolutionary army in Hunan, Mao spent six months as a soldier. While he probably had not yet clearly grasped the idea that, as he later put it, "political power grows out of the barrel of a gun," his first brief military experience at least confirmed his boyhood admiration of military leaders and exploits. In primary school days, his heroes had included not only the great warrior-emperors of the Chinese past but Napoleon and George Washington as well.

The spring of 1912 saw the birth of the new Chinese Republic and the end of Mao's military service. For a year he drifted from one thing to another, trying, in turn,

a police school, a law school, and a business school; he studied history in a secondary school and then spent some months reading many of the classic works of the Western liberal tradition in the provincial library. This period of groping, rather than indicating any lack of decision in Mao's character, was a reflection of China's situation at the time. The abolition of the official civil service examination system in 1905 and the piecemeal introduction of Western learning in so-called modern schools had left young people in a state of uncertainty as to what type of training, Chinese or Western, could best prepare them for a career or for service to their country.

Mao eventually was graduated from the First Provincial Normal School in Ch'ang-sha in 1918. While officially an institution of secondary level rather than of higher education, the normal school offered a high standard of instruction in Chinese history, literature, and philosophy as well as in Western ideas. While at the school, Mao also acquired his first experience in political activity by helping to establish several student organizations. The most important of these was the New People's Study Society, founded in the winter of 1917-18, many of whose members were later to join the Communist Party.

From the normal school in Ch'ang-sha, Mao went to Peking University, China's leading intellectual centre. The half year he spent in Peking working as a librarian's assistant was of disproportionate importance in shaping his future career, for it was then that he came under the influence of the two men who were to be the principal figures in the foundation of the Chinese Communist Party: Li Dazhao (Li Ta-chao) and Chen Duxiu (Ch'en Tu-hsiu). Moreover, he found himself at Peking University precisely during the months leading up to the May Fourth Movement of 1919, which was to a considerable extent the fountainhead of all of the changes that were to take place in China in the ensuing half century.

In a limited sense, May Fourth Movement is the name given to the student demonstrations protesting against the Paris Peace Conference's decision to hand over former German concessions in Shantung Province to Japan instead of returning them to China. But the term also evokes a period of rapid political and cultural change, beginning in 1915, that resulted in the Chinese radicals' abandonment of Western liberalism for Marxism-Leninism as the answer to China's problems and the subsequent founding of the Chinese Communist Party in 1921. The shift from the difficult and esoteric classical written language to a far more accessible vehicle of literary expression patterned on colloquial speech also took place during this period. At the same time, a new and very young generation moved to the centre of the political stage. To be sure, the demonstration on May 4 was launched by Chen Duxiu, but the students soon realized that they themselves were the main actors. In an editorial published in July 1919, Mao wrote:

The world is ours, the nation is ours, society is ours. If we do not speak, who will speak? If we do not act, who will act?

From then onward, his generation never ceased to regard itself as responsible for the nation's fate, and, indeed, its members remained in power, both in Peking and in Taipei, until the 1970s.

During the summer of 1919 Mao Zedong helped to establish in Ch'ang-sha a variety of organizations that brought the students together with the merchants and the workers—but not yet with the peasants—in demonstrations aimed at forcing the government to oppose Japan. His writings at the time are filled with references to the "army of the red flag" throughout the world and to the victory of the Russian Revolution, but it was not until January 1921 that he was finally committed to Marxism as the philosophical basis of the revolution.

Mao and the Chinese Communist Party. In September

May
Fourth
Movement

Enlistment
in revolu-
tionary
army

1920 he became principal of the Lin Ch'ang-sha primary school, and in October he organized a branch of the Socialist Youth League there. That winter he married Yang Kaihui (Yang K'ai-hui), the daughter of his former ethics teacher. In July 1921 he attended the First Congress of the Chinese Communist Party, together with representatives from the other Communist groups in China and two delegates from the Moscow-based Comintern (Communist International). In 1923, when the young party entered into an alliance with Sun Yat-sen's Kuomintang (Nationalist Party), Mao was one of the first Communists to join the Kuomintang and to work within it. During the first half of 1924, he lived mostly with his wife and two infant sons in Shanghai, where he was a leading member of the Kuomintang Executive Bureau.

In the winter of 1924-25, Mao returned to his native village of Shao-shan for a rest. There, after witnessing demonstrations by peasants stirred into political consciousness by the shooting of several dozen Chinese by foreign police in Shanghai (May and June 1925), Mao suddenly became aware of the revolutionary potential inherent in the peasantry. Although born in a peasant household, he had, in the course of his student years, adopted the Chinese intellectual's traditional view of the workers and peasants as ignorant and dirty. His conversion to Marxism had forced him to revise his estimate of the urban proletariat, but he continued to share Marx's own contempt for the backward and amorphous peasantry. Now he turned back to the rural world of his youth as the source of China's regeneration. Following the example of other Communists working within the Kuomintang who had already begun to organize the peasants, Mao sought to channel the spontaneous protests of the Hunanese peasants into a network of peasant associations.

The Communists and the Kuomintang. Pursued by the military governor of Hunan, Mao was soon forced to flee his native province once more, and he returned for another year to an urban environment—this time to Canton, the main power base of the Kuomintang. But, though he lived in Canton, Mao still focused his attention on the countryside. He became the acting head of the propaganda department of the Kuomintang—in which capacity he edited its leading organ, the *Political Weekly*, and attended the Second Kuomintang Congress in January 1926—but he also served at the Peasant Movement Training Institute, set up in Canton under the auspices of the Kuomintang, as principal of the sixth training session. Chiang Kai-shek had become the leader of the Kuomintang after the death of Sun Yat-sen in March 1925; and although Chiang still declared his allegiance to the "world revolution" and wished to avail himself of Soviet aid, he was determined to remain master in his own house. He therefore expelled most Communists from responsible posts in the Kuomintang in May 1926. Mao, however, stayed on at the institute until October of that year. Most of the young peasant activists Mao trained were shortly at work strengthening the position of the Communists.

In July 1926, Chiang Kai-shek set out on what became known as the Northern Expedition, aiming to unify the country under his own leadership and to overthrow the conservative government in Peking as well as other warlords. In November Mao once more returned to Hunan; there, in January and February 1927, he investigated the peasant movement and concluded that in a very short time several hundred million peasants in China would "rise like a tornado or tempest—a force so extraordinarily swift and violent that no power, however great, will be able to suppress it." Strictly speaking, this prediction proved to be false. Revolution in the shape of spontaneous action by hundreds of millions of peasants did not sweep across China "in a very short time," or indeed at all. Chiang Kai-shek, who was bent on an alliance with the propertied classes in the cities and in the countryside, turned against the worker and peasant revolution, and in April he massacred the very Shanghai workers who had delivered the city to him. Stalin's strategy for carrying out revolution in alliance with the Kuomintang collapsed, and the Chinese Communist Party was virtually annihilated in the cities and decimated in the countryside. But in a broader

and less literal sense, Mao's prophecy was justified. In October 1927 Mao led a few hundred peasants who had survived the autumn harvest uprising in Hunan to a base in the Ching-kang Shan (Ching-kang Mountains), on the Kiangsi-Hunan border, and embarked on a new type of revolutionary warfare in the countryside in which the Red Army, rather than the unarmed masses, would play the central role. But it was only because a large proportion of China's hundreds of millions of peasants sympathized with and supported this effort that Mao Zedong was able in the course of the civil war to encircle the cities from the countryside and thus defeat Chiang Kai-shek and gain control of the country.

The road to power. Mao Zedong's 22 years in the wilderness can be divided into four phases. The first of these is the initial three years when Mao and Zhu De (Chu Teh), the commander in chief of the army, successfully developed the tactics of guerrilla warfare from base areas in the countryside. These activities, however, were regarded even by their protagonists, and still more by the Central Committee in Shanghai (and by the Comintern in Moscow), as a holding operation until the next upsurge of revolution in the urban centres. In the summer of 1930 the Red Army was ordered by the Central Committee to occupy several major cities in south central China in the hope of sparking a revolution by the workers. When it became evident that persistence in this attempt could only lead to further costly losses, Mao disobeyed orders and abandoned the battle to return to the base in southern Kiangsi. During this year Mao's wife was executed by the Kuomintang and he married He Zizhen (Ho Tzu-chen), with whom he had been living since 1928.

The second phase (the Kiangsi period) centres on the founding of the Chinese Soviet Republic, November 1931, in a portion of Kiangsi Province, with Mao as chairman. Since there was little support for the revolution in the cities, the promise of ultimate victory now seemed to reside in the gradual strengthening and expansion of the base areas. The Soviet regime soon came to control a population of several million; the Red Army, grown to a strength of some 200,000, easily defeated large forces of inferior troops sent against it by Chiang Kai-shek in the first four of the so-called encirclement and annihilation campaigns. But it was unable to stand up against Chiang's own elite units, and in October 1934 the major part of the Red Army, Mao, and his pregnant wife abandoned the base in Kiangsi and set out for the northwest of China, on what is known as the Long March.

There is wide disagreement among specialists as to the extent of Mao's real power, especially in the years 1932-34, and as to which military strategies were his or other party leaders'. The majority view is that, in the last years of the Chinese Soviet Republic, Mao functioned to a considerable extent as a figurehead with little control over policy, especially in military matters. In any case, he achieved de facto leadership over the party (though not the formal title of chairman) only at the Tsun-i Conference of January 1935 during the Long March.

When some 8,000 troops who had survived the perils of the Long March arrived in Shensi Province in northwestern China in the autumn of 1935, events were already moving toward the third phase in Mao's rural odyssey, which was to be characterized by a renewed united front with the Kuomintang against Japan and by the rise of Mao to unchallenged supremacy in the party. This phase is often called the Yen-an period (for the town in Shensi where the Communists were based), although Mao did not move to Yen-an until December 1936. In August 1935 the Comintern at its Seventh Congress in Moscow proclaimed the principle of an anti-Fascist united front, and in May 1936 the Chinese Communists for the first time accepted the prospect that such a united front might include Chiang Kai-shek himself, and not merely dissident elements in the Nationalist camp. The so-called Sian Incident of December 1936, in which Chiang was kidnapped by military leaders from northeastern China who wanted to fight Japan and recover their homelands rather than participate in civil war against the Communists, accelerated the evolution toward unity. By the time the Japanese

Struggle
with
Chiang
Kai-shek

Peasants
and the
revolution

The
Kiangsi
period

Yen-an
period

began their attempt to subjugate all of China in July 1937, the terms of a new united front between the Communists and the Kuomintang had been virtually settled, and the formal agreement was announced in September 1937.

In the course of the anti-Japanese war, the Communists broke up a substantial portion of their army into small units and sent them behind the enemy lines to serve as nuclei for guerrilla forces that effectively controlled vast areas of the countryside, stretching between the cities and communication lines occupied by the invader. As a result, they not only expanded their military forces to somewhere between 500,000 and 1,000,000 at the time of the Japanese surrender but also established effective grassroots political control over a population that may have totaled as many as 90,000,000. It has been argued that the support of the rural population was won purely by appeals to their nationalist feeling in opposition to the Japanese. This certainly was fundamental, but Communist agrarian policies likewise played a part in securing broad support among the peasantry.

Writings,
1936–40

During the years 1936–40, Mao had, for the first time since the 1920s, the leisure to devote himself to reflection and writing. It was then that he first read in translation a certain number of Soviet writings on philosophy and produced his own account of dialectical materialism, of which the best known portions are those entitled "On Practice" and "On Contradiction." More important, Mao produced the major works that synthesized his own experience of revolutionary struggle and his vision of how the revolution should be carried forward in the context of the united front. On military matters there was first *Strategic Problems of China's Revolutionary War*, written in December 1936 to sum up the lessons of the Kiangsi period (and also to justify the correctness of his own military line at the time), and then *On Protracted War* and other writings of 1938 on the tactics of the anti-Japanese war. As to his overall view of the events of these years, Mao adopted an extremely conciliatory attitude toward the Kuomintang in his report entitled *On the New Stage* (October 1938), in which he attributed to it the leading role both in the war against Japan and in the ensuing phase of national reconstruction. By the winter of 1939–40, however, the situation had changed sufficiently so that he could adopt a much firmer line, claiming leadership for the Communists. Internationally, Mao argued, the Chinese revolution was a part of the world proletarian revolution directed against imperialism (whether it be British, German, or Japanese); internally, the country should be ruled by a "joint dictatorship of several parties" belonging to the anti-Japanese united front. For the time being, Mao felt, the aims of the Communist Party coincided with the aims of the Kuomintang, and therefore Communists should not try to rush ahead to socialism and thus disrupt the united front. But neither should they have any doubts about the ultimate need to take power into their own hands in order to move forward to socialism. During this period, in 1939, Mao divorced He Zizhen and married a well-known film actress, Lan P'ing, later called Jiang Qing (Chiang Ch'ing).

The issues of Kuomintang–Communist rivalry for the leadership of the united front are related to the continuing struggle for supremacy within the Chinese Communist Party, for Mao's two chief rivals—Wang Ming, who had just returned from a long stay in Moscow, and Chang Kuo-t'ao, who had at first refused to accept Mao's political and military leadership—were both accused of excessive slavishness toward the Kuomintang. But perhaps even more central in Mao's ultimate emergence as the acknowledged leader of the party was the question of what he had called in October 1938 the "Sinification" of Marxism—its adaptation not only to Chinese conditions but to the mentality and cultural traditions of the Chinese people.

Mao could not claim the firsthand knowledge possessed by many other leading members of the Chinese Communist Party of how Communism worked within the Soviet Union nor the ability to read Marx or Lenin in the original, which some of them enjoyed. He could and did claim, however, to know and understand China. The differences between him and the Soviet-oriented faction in the party came to a head at the time of the so-called

Rectification Campaign of 1942–43. This program aimed at giving a basic grounding in Marxist theory and Leninist principles of party organization to the many thousands of new members who had been drawn into the party in the course of the expansion since 1937. But a second and equally important aspect of the movement was the elimination of what Mao called "foreign dogmatism"—in other words, blind imitation of Soviet experience and obedience to Soviet directives.

In March 1943 Mao achieved for the first time formal supremacy over the party, becoming chairman of the Secretariat and of the Politburo. Shortly thereafter, the Rectification Campaign took, for a time, the form of a harsh purge of elements not sufficiently loyal to Mao. The campaign was run by Kang Sheng (K'ang Sheng), who was later to be one of Mao's key supporters in the Cultural Revolution. Exaggerating considerably this dimension of events, Soviet spokesmen have bitterly denounced the Rectification Campaign as an attempt to purge the Chinese Communist Party of all those elements genuinely imbued with "proletarian internationalism" (*i.e.*, devotion to Moscow). It is therefore not surprising that, as Mao's campaign in the countryside moved into its fourth and last phase—that of civil war with the Kuomintang—Stalin's lack of enthusiasm for a Chinese Communist victory should have become increasingly evident. Looking back at this period in 1962, when the Sino-Soviet conflict had come to a head, Mao declared:

In 1945, Stalin wanted to prevent China from making revolution, saying that we should not have a civil war and should cooperate with Chiang Kai-shek, otherwise the Chinese nation would perish. But we did not do what he said. The revolution was victorious. After the victory of the revolution he [Stalin] next suspected China of being a Yugoslavia, and that I would become a second Tito.

This account of Stalin's attitude is substantiated by a whole series of public gestures at the time, culminating in the fact that when the People's Liberation Army took Nanking in April 1949, the Soviet ambassador was the only foreign diplomat to accompany the retreating Nationalist government to Canton. Stalin's motives were obviously those described by Mao in the above passage; he did not believe in the capacity of the Chinese Communists to achieve a clear-cut victory, and he thought they would be a nuisance if they did.

Formation of the People's Republic of China. Nevertheless, when the Communists did take power in China, both Mao and Stalin had to make the best of the situation. In December 1949 Mao, now chairman of the People's Republic of China—which he had proclaimed on October 1—traveled to Moscow, where, after two months of arduous negotiations, he succeeded in persuading Stalin to sign a treaty of mutual assistance accompanied by limited economic aid. Before the Chinese had time to profit from the resources made available for economic development, however, they found themselves dragged into the Korean War in support of the Moscow-oriented regime in P'yongyang. Only after this baptism of fire did Stalin, according to Mao, begin to have confidence in him and believe he was not first and foremost a Chinese nationalist.

Despite these tensions with Moscow, the policies of the People's Republic of China in its early years were based in very many respects, as Mao later said, on "copying from the Soviets." While Mao and his comrades had experience in guerrilla warfare, in mobilization of the peasants in the countryside, and in political administration at the grass roots, they had no firsthand knowledge of running a state or of large-scale economic development. In such circumstances, the Soviet Union provided the only available model. A five-year plan was therefore drawn up under Soviet guidance; it was put into effect in 1953 and included Soviet technical assistance and a number of complete industrial plants. Yet, within two years, Mao had taken steps that were to lead to the breakdown of the political and ideological alliance with Moscow.

The emergence of Mao's road to socialism. In the spring of 1949, Mao proclaimed that while in the past the Chinese revolution had followed the unorthodox path of "encircling the cities from the countryside," it would

Mao
becomes
chairman

in the future take the orthodox road of the cities leading and guiding the countryside. In harmony with this view, he had agreed in 1950 with Liu Shaoqi (Liu Shao-ch'i) that collectivization would be possible only when China's heavy industry had provided the necessary equipment for mechanization. In a report of July 1955, he reversed this position, arguing that in China the social transformation could run ahead of the technical transformation. Deeply impressed by the achievements of certain cooperatives that claimed to have radically improved their material conditions without any outside assistance, he came to believe in the limitless capacity of the Chinese people, especially of the rural masses, to transform at will both nature and their own social relations when mobilized for revolutionary goals. Those in the leadership who did not share this vision he denounced as "old women with bound feet." He made these criticisms before an ad hoc gathering of provincial and local party secretaries, thus creating a ground swell of enthusiasm for rapid collectivization such that all those in the leadership who had expressed doubts about Mao's ideas were soon presented with a fait accompli. The tendency thus manifested to pursue his own ends outside the collective decision-making processes of the party was to continue and to be accentuated.

Even before Nikita S. Khrushchev's secret speech of February 1956 denouncing Stalin's crimes, Mao Zedong and his colleagues had been discussing measures for improving the morale of the intellectuals in order to secure their willing participation in building a new China. At the end of April, Mao proclaimed the policy of "letting a hundred flowers bloom"—that is, the freedom to express many diverse ideas—designed to prevent the development in China of a repressive political climate analogous to that in the Soviet Union under Stalin. In the face of the disorders called forth by de-Stalinization in Poland and Hungary, Mao did not retreat but rather pressed boldly forward with this policy, against the advice of many of his senior colleagues, in the belief that the contradictions that still existed in Chinese society were mainly nonantagonistic. When the resulting "great blooming and contending" got out of hand and called into question the axiom of party rule, Mao savagely turned against the educated elite, which he felt had betrayed his confidence. Henceforth, he would rely primarily on the creativity of the rank and file as the agent of modernization. As for the specialists, if they were not yet sufficiently "red," he would remold them by sending them to work in the countryside.

It was against this background that Mao, during the winter of 1957–58, worked out the policies that were to characterize the Great Leap Forward, formally launched in May 1958. While his economic strategy was by no means so one-sided and simplistic as commonly believed in the 1960s and '70s, and though he still proclaimed industrialization and a "technical revolution" as his goals, Mao displayed continuing anxiety regarding the corrupting influence of the fruits of technical progress and an acute nostalgia for the purity and egalitarianism that had marked the moral and political world of the Ching-kang Shan and Yen-an eras.

Thus it was logical that he should endorse and promote the establishment of "people's communes" as part of the Great Leap strategy. As a result, the peasants, who had been organized into cooperatives in 1955–56 and then into fully socialist collectives in 1956–57, found their world turned upside down once again in 1958. Neither the resources nor the administrative experience necessary to operate such enormous new social units of several thousand households were in fact available, and not surprisingly the consequences of these changes were chaos and economic disaster.

By the winter of 1958–59, Mao himself had come to recognize that some adjustments were necessary, including decentralization of ownership to the constituent elements of the communes and a scaling down of the unrealistically high production targets in both industry and agriculture. He insisted, however, that in broad outline his new Chinese road to socialism, including the concept of the communes and the belief that China, though "poor and blank," could leap ahead of other countries, was basically

sound. At the Lushan meeting of the Central Committee in July–August 1959, Peng Dehuai (P'eng Te-huai), the minister of defense, denounced the excesses of the Great Leap and the economic losses they had caused. He was immediately removed from all his party and state posts and remained in detention until his death during the Cultural Revolution. From that time, Mao regarded any criticism of his policies as nothing less than a crime of *lèse-majesté*, meriting exemplary punishment.

Retreat and counterattack. Though few spoke up at Lushan in support of Peng, a considerable number of the top leaders sympathized with him in private. Almost immediately, in 1960, Mao began building an alternative power base in the People's Liberation Army, which the new defense minister, Lin Biao (Lin Piao), had set out to turn into a "great school of Mao Zedong Thought." At about the same time, Mao began to denounce the emergence, not only in the Soviet Union but also in China itself, of "new bourgeois elements" among the privileged strata of the state and party bureaucracy and the technical and artistic elite. Under these conditions, he concluded, a "protracted, complex, and sometimes even violent class struggle" would continue during the whole socialist stage.

The open split with the Soviet Union, which had become public and irreparable by 1963—though it can be traced to Mao's resentment at Khrushchev's failure to consult him before launching de-Stalinization—resulted, above all, from the Soviet reaction to the Great Leap policies. Regarding Mao's claims for the communes as ideologically presumptuous, Khrushchev heaped ridicule upon them; he underlined his displeasure by withdrawing Soviet technical assistance in 1960, leaving many large plants unfinished. Khrushchev also tried to put pressure on China in its dealings with Taiwan and India and in other foreign-policy issues. Mao forgot neither the affront to his and China's dignity nor the economic damage.

As for class struggle in China itself, Mao's fear that revisionism might appear there was also heightened by the policies pursued in the early 1960s to deal with the economic consequences of the Great Leap Forward. The disorganization and waste created by the Great Leap, compounded by natural disasters and by the termination of Soviet economic aid, led to widespread famine in which, according to much later official Chinese accounts, millions of people died. The response to this situation by Liu Shaoqi (who had succeeded Mao as chairman of the People's Republic in 1959), Deng Xiaoping (Teng Hsiao-p'ing), and the economic planners was to make use of material incentives and to strengthen the role of individual households in agricultural production. At first Mao agreed reluctantly that such steps were necessary, but during the first half of 1962 he came increasingly to perceive the methods used to promote recovery as implying the repudiation of the whole thrust of the Great Leap strategy. It was as a direct response to this challenge that at the 10th Plenary Session of the Central Committee in September 1962 he issued the call, "Never forget the class struggle!"

During the next three years Mao waged such a struggle, primarily through the Socialist Education Movement in the countryside, and it was over the guidelines for this campaign that the major political battles were fought within the Chinese leadership. At the end of 1964, when Liu Shaoqi refused to accept Mao's demand to direct the main thrust of class struggle against "capitalist roaders" in the party, Mao decided that "Liu had to go."

The Cultural Revolution. The movement that became known as the Great Proletarian Cultural Revolution represented an attempt by Mao to go beyond the party rectification campaigns, of which there had been many since 1942, and to devise a new and more radical method for dealing with what he saw as the bureaucratic degeneration of the party. But it also represented, beyond any doubt or question, a deliberate effort to eliminate those in the leadership who, over the years, had dared to cross him. The victims, from throughout the party hierarchy, suffered more than mere political disgrace. All were publicly humiliated and detained for varying periods, sometimes under very harsh conditions; many were beaten and tortured, and not a few were killed or driven to suicide.

Hundred
Flowers
speech

Great Leap
Forward

Split with
Soviet
Union

Elimination of Liu Among the casualties was Liu, who died because he was denied proper medical attention.

The justification for these sacrifices was defined in a key slogan of the time: "Fight selfishness, criticize revisionism." When the Red Guards, who constituted the first shock troops of Mao's enterprise, burst upon the scene in the summer of 1966, with their battle cry, "To rebel is justified!" it seemed for a time that not only the power of the party cadres but also authority in all its forms was being questioned. It soon became evident that Mao, who in 1956 had justified decentralization as a means to building a "strong socialist state," still believed in the need for state power. When the Shanghai leftists Zhang Chunqiao (Chang Ch'un-ch'iao) and Yao Wenyuan (Yao Wen-yüan)—who were later to make up half the Gang of Four—came to see him in February 1967, immediately after setting up the Shanghai Commune, Mao asserted that the demand for the abolition of "heads" (leaders), which had been heard in their city, was "extreme anarchism" and "most reactionary"; in fact, he stated, there would "always be heads." Communes, he added, were "too weak when it came to suppressing counterrevolution" and in any case required party leadership. He therefore ordered them to dissolve theirs and to replace it with a "revolutionary committee."

These committees, based on an alliance of former party cadres, young activists, and representatives of the People's Liberation Army, were to remain in place until two years after Mao's death. At first, they were largely controlled by the army. The Ninth Congress of 1969 initiated the process of rebuilding the party; and the death of Lin Biao diminished, though it by no means eliminated, the army's role. Thereafter, it seemed briefly, in 1971–72, that a compromise, of which Zhou Enlai (Chou En-lai) was the architect, might produce some kind of synthesis between the values of the Cultural Revolution and the pre-1966 political and economic order.

Even before Zhou's death in January 1976, however, this compromise had been overturned. All recognition by Mao of the importance of professional skills was swallowed up in an orgy of political rhetoric, and all things foreign were regarded as counterrevolutionary. Mao's last decade, which had opened with manifestos in favour of the Paris Commune model of mass democracy, closed with paeans of praise to that most implacable of centralizing despots, Shih Huang-ti, the first Ch'in emperor. Mao Zedong died in Peking on Sept. 9, 1976.

Assessment. While the Cultural Revolution was an entirely logical culmination of Mao's last two decades, it was by no means the only possible outcome of his approach to revolution, nor need a judgment of his work as a whole be based primarily on this last phase.

Few would deny Mao Zedong the major share of credit for devising the pattern of struggle based on guerrilla warfare in the countryside that ultimately led to victory in the civil war and thereby to the overthrow of the Kuomintang, the distribution of land to the peasants, and the restoration of China's independence and sovereignty. These achievements must be given a weight commensurate with the degree of injustice prevailing in Chinese society before the revolution and with the humiliation felt by the Chinese people as a result of the dismemberment of their country

by the foreign powers. "We have stood up," Mao said in September 1949. These words will not be forgotten.

Mao's record after 1949 is more ambiguous. The official Chinese view, defined in June 1981, is that his leadership was basically correct until the summer of 1957, but from then on it was mixed at best and frequently quite wrong. It cannot be disputed that Mao's two major innovations of his later years, the Great Leap and the Cultural Revolution, were ill-conceived and led to disastrous consequences. His goals of combating bureaucracy, encouraging popular participation, and stressing China's self-reliance were generally laudable, but the methods he used to pursue them, though bold and imaginative, were largely self-defeating.

Looking at Mao's whole career, it is not easy to put a figure on the positive and negative aspects. How does one weigh the good fortune of peasants acquiring land against millions of executions and deaths from civil war? How does one balance the real economic achievements after 1949 against the starvation that came in the wake of the Great Leap Forward or the bloody shambles of the Cultural Revolution? It is, perhaps, possible to accept the official verdict that, despite the "errors of his later years," Mao's merits outweighed his faults, while underscoring the fact that the account is very finely balanced.

BIBLIOGRAPHY. A detailed biography of Mao Zedong is ROSS TERRILL, *Mao* (1980, reissued 1981). Two earlier works that remain useful for the pre-1949 period are JEROME CH'EN, *Mao and the Chinese Revolution* (1965, reissued 1972); and STUART R. SCHRAM, *Mao Tse-tung*, rev. ed. (1967, reprinted 1974). See also Schram's *Mao Zedong, A Preliminary Reassessment* (1983). The most vivid account of Mao's youth is his autobiography as recounted in 1936 in EDGAR SNOW, *Red Star over China*, rev. ed. (1968, reissued 1974). JUI LI, *The Early Revolutionary Activities of Comrade Mao Tse-tung* (1977; originally published in Chinese, 1957), is also an important source.

Regarding Mao Zedong's thought, a substantial collection of source materials for the period before 1949 is available in *Selected Works of Mao Tse-tung*, 5 vol. (1961–77). A variorum in Chinese of the collected writings of Mao to 1949 is *Mao Tse-tung chi*, ed. by MINORU TAKEUCHI, 10 vol. (1970–72), completed by a set of supplements, *Mao Tse-tung chi pu chüan* (1983–85). Mao's talks and letters from 1956 to 1971 are found in STUART R. SCHRAM (ed.), *Mao Tse-tung Unrehearsed* (1974, U.S. title, *Chairman Mao Talks to the People*, 1975). See also JEROME CH'EN (ed.), *Mao Papers* (1970); and MAO TSETUNG, *A Critique of Soviet Economics* (1977), trans. from Chinese by MOSS ROBERTS. On Mao's thought in his early years, see BRANTLY WOMACK, *The Foundation of Mao Zedong's Political Thought, 1917–1935* (1982); FREDERIC WAKEMAN, JR., *History and Will: Philosophical Perspectives of Mao Tse-tung's Thought* (1973, reprinted 1975), which links Mao's ideas of the May Fourth period with those of the Cultural Revolution; and RAYMOND F. WYLIE, *The Emergence of Maoism* (1980). JOHN BRYAN STARR, *Continuing the Revolution: The Political Thought of Mao* (1979), is a comprehensive overview that accepts at face value the Chinese view of the Chairman during his lifetime. Among the older works, ARTHUR A. COHEN, *The Communism of Mao Tse-tung* (1964, reprinted 1971), stresses the Stalinist roots of Mao's thought; and JAMES HSIUNG, *Ideology and Practice: The Evolution of Chinese Communism* (1970), emphasizes the links between Mao's thought and Chinese tradition. See also STUART R. SCHRAM (ed.), *The Political Thought of Mao Tse-tung*, rev. ed. (1969). Finally, a series of useful if somewhat premature appreciations are in DICK WILSON (ed.), *Mao Tse-tung in the Scales of History* (1977).

(S.R.S.)

Mao's death

Mapping and Surveying

Maps and charts are representations of features on the Earth's surface drawn to scale. Globes are maps presented on the surface of a sphere. The definition can be extended to include representations of the Moon, Mars, and other bodies. Star charts and related celestial plats, long used for navigation and astronomical studies, might also be included; they are treated in the article STARS AND STAR CLUSTERS.

In order to imply the elements of accurate relationships, and some formal method of projecting the spherical subject to a map plane, further qualifications might be applied to the definition. The tedious and somewhat abstract statements resulting from attempts to formulate precise definitions of maps and charts are more likely to confuse than to clarify. The words map, chart, and plat are used somewhat interchangeably. The connotations of use, however, are distinctive: charts for navigation purposes (nautical and aeronautical), plats (in a property-boundary sense) for land-line references and ownership, and maps for general reference.

Cartography is the art and science of making maps and charts. As such, it is allied with geography in its concern with the broader aspects of the Earth and its life. In early times cartographic efforts were more artistic than scientific and factual. As man explored and recorded his environment, the quality of his maps and charts improved. These lines of Jonathan Swift were inspired by early maps:

So geographers, in Afric maps,
With savage pictures fill their gaps,
And o'er uninhabitable downs
Place elephants for want of towns.

Topographic maps are graphic representations of natural and man-made features of parts of the Earth's surface plotted to scale. They show the shape of land and record elevations above sea level, lakes, streams and other hydrographic features, and roads and other works of man. In short, they provide a complete inventory of the terrain and important information for all activities involving the use and development of the land. They provide the bases for specialized maps and data for compilation of generalized maps of smaller scale.

Nautical charts are maps of coastal and marine areas, providing information for navigation. They include depth curves or soundings or both; aids to navigation such as buoys, channel markers, and lights; islands, rocks, wrecks, reefs and other hazards; and significant features of the coastal areas, including promontories, church steeples, water towers, and other features helpful in determining positions from offshore.

The terms hydrography and hydrographer date from the mid-16th century; their focus has become restricted to studies of ocean depths and of the directions and intensities of oceanic currents; though at various times they embraced much of the sciences now called hydrology and oceanography. The British East India Company employed hydrographers in the 18th century, and the first hydrographer of the Royal Navy, Alexander Dalrymple (1737–1808), was appointed in 1795. A naval observatory and hydrographic office was established administratively in the United States Navy in 1854. In 1866 a hydrographic office was established by statute, and in 1962 it was renamed the U.S. Naval Oceanographic Office.

Interest in the charting of oceanic areas away from sea-coasts developed in the second half of the 19th century, concurrently with the perfection of submarine cables. As knowledge of the configuration of the ocean basins increased, the attention of scientists was drawn to this field of study. A feature of marine science since the 1950s has been increasingly detailed bathymetric (water-depth measurement) surveys of selected portions of the seafloor. Together with collection of associated geophysical data and sampling of sediments, these studies assist in interpreting the geologic history of the ocean-covered portion of the Earth's crust.

Aeronautical charts provide essential data for the pilot and air navigator. They are, in effect, small-scale topographic maps on which current information on aids to navigation have been superimposed. To facilitate rapid recognition and orientation, principal features of the land that would be visible from an aircraft in flight are shown to the exclusion of less important details.

This article is organized into the following sections:

History of cartography 473	Types and uses of maps and charts 482
Maps and geography in the ancient world 474	World status of mapping and basic data
Greek maps and geography	Types of maps and charts available
The Roman period	Government and other mapping agencies
The Middle Ages 475	International organizations
The age of discovery and exploration 476	Modern mapmaking techniques 484
Revival of Ptolemy	Compilation from existing materials
Maps of the discoveries	Map production from original surveys
18th century to the present 477	Final steps in map preparation
The rise of national surveys	Automation in mapping
International Map of the World (IMW)	Surveying 486
World War II and after	History 487
Mapmaking 478	Modern surveying 488
Elements 478	Basic control surveys
Map scales and classifications	Global positioning
Map projections	Establishing the framework
Development of reference spheroids	Detail surveying
Geographic and plane coordinate systems	Aerial surveying
Basic data for compilation	Hydrography
Symbolization	Height determination
Nomenclature	Bibliography 494

History of cartography

Centuries before the Christian Era, Babylonians drew maps on clay tablets, of which the oldest specimens found so far have been dated about 2300 BC. This is the earliest positive evidence of graphic representations of parts of

the Earth; it may be assumed that mapmaking goes back much further and that it began among nonliterate peoples. It is logical to assume that men very early made efforts to communicate with each other regarding their environment by scratching routes, locations, and hazards on the ground and later on bark and skins.

The earliest maps must have been based on personal experience and familiarity with local features. They doubtless showed routes to neighbouring tribes, where water and other necessities might be found, and the locations of enemies and other dangers. Nomadic life stimulated such efforts by recording ways to cross deserts and mountains, the relative locations of summer and winter pastures, and dependable springs, wells, and other information.

Markings on cave walls that are associated with paintings by primitive man have been identified by some archaeologists as attempts to show the game trails of the animals depicted, though there is no general agreement on this. Similarly, networks of lines scratched on certain bone tablets could possibly represent hunting trails, but there is definitely no conclusive evidence that the tablets are indeed maps.

Many nonliterate peoples, however, are skilled in depicting essential features of their localities and travels. During Capt. Charles Wilkes's exploration of the South Seas in the 1840s, a friendly islander drew a good sketch of the whole Tuamotu Archipelago on the deck of the captain's bridge. In North America the Pawnee Indians were reputed to have used star charts painted on elk skin to guide them on night marches across the plains. Montezuma is said to have given Cortés a map of the whole Mexican Gulf area painted on cloth, while Pedro de Gamboa reported that the Incas used sketch maps and cut some in stone to show relief features. Many specimens of early Eskimo sketch maps on skin, wood, and bone have been found.

MAPS AND GEOGRAPHY IN THE ANCIENT WORLD

The earliest specimens thus far discovered that are indisputably portrayals of land features are the Babylonian tablets previously mentioned; certain land drawings found in Egypt and paintings discovered in early tombs are nearly as old. It is quite probable that these two civilizations developed their mapping skills more or less concurrently and in similar directions. Both were vitally concerned with the fertile areas of their river valleys and therefore doubtless made surveys and plats soon after settled communities were established. Later they made plats for the construction of canals, roads, and temples—the equivalent of today's engineering plans.

A tablet unearthed in Iraq shows the Earth as a disk surrounded by water with Babylon as its centre. Aside from this specimen, dating from about 1000 bc, there appear to have been rather few attempts by Babylonians and Egyptians to show the form and extent of the Earth as a whole. Their mapmaking was preoccupied with more practical needs, such as the establishment of boundaries. Not until the time of the Greek philosopher-geographers did speculations and conclusions as to the nature of the Earth begin to take form.

Greek maps and geography. The Greeks were outstanding among peoples of the ancient world for their pursuit and development of geographic knowledge. The shortage of arable land in their own region led to maritime explo-

ration and the development of commerce and colonies. By 600 bc Miletus, on the Aegean, had become a centre of geographic knowledge, as well as of cosmographic speculation.

Hecataeus, a scholar of Miletus, probably produced the first book on geography in about 500 bc. A generation later Herodotus, from more extensive studies and wider travels, expanded upon it. A historian with geographic leanings, Herodotus recorded, among other things, an early circumnavigation of the African continent by Phoenicians. He also improved on the delineation of the shape and extent of the then-known regions of the world, and he declared the Caspian to be an inland sea, opposing the prevailing view that it was part of the "northern oceans" (Figure 1).

Herodotus

Although Hecataeus regarded the Earth as a flat disk surrounded by ocean, Herodotus and his followers questioned the concept and proposed a number of other possible forms. Indeed, the philosophers and scholars of the time appear to have been preoccupied for a number of years with discussions on the nature and extent of the world. Some modern scholars attribute the first hypothesis of a spherical Earth to Pythagoras (6th century bc) or Parmenides (5th century). The idea gradually developed into a consensus over many years. In any case by the mid-4th century the theory of a spherical Earth was well accepted among Greek scholars, and about 350 bc Aristotle formulated six arguments to prove that the Earth was, in truth, a sphere. From that time forward, the idea of a spherical Earth was generally accepted among geographers and other men of science.

About 300 bc Dicaearchus, a disciple of Aristotle, placed an orientation line on the world map, running east and west through Gibraltar and Rhodes. Eratosthenes, Marinus of Tyre, and Ptolemy successively developed the reference-line principle until a reasonably comprehensive system of parallels and meridians, as well as methods of projecting them, had been achieved.

The greatest figure of the ancient world in the advancement of geography and cartography was Claudius Ptolemaeus (Ptolemy; AD 90–168). An astronomer and mathematician, he spent many years studying at the library in Alexandria, the greatest repository of scientific knowledge at that time. His monumental work, the *Guide to Geography* (*Geōgraphikē hyphēgēsis*), was produced in eight volumes. The first volume discussed basic principles and dealt with map projection and globe construction. The next six volumes carried a list of the names of some 8,000 places and their approximate latitudes and longitudes. Except for a few that were made by observations, the greater number of these locations were determined from older maps, with approximations of distances and directions taken from travelers. They were accurate enough to show relative locations on the very small-scale, rudimentary maps that existed.

The eighth volume was a most important contribution, containing instructions for preparing maps of the world and discussions on mathematical geography and other

By courtesy of the Library of Congress, Washington D C

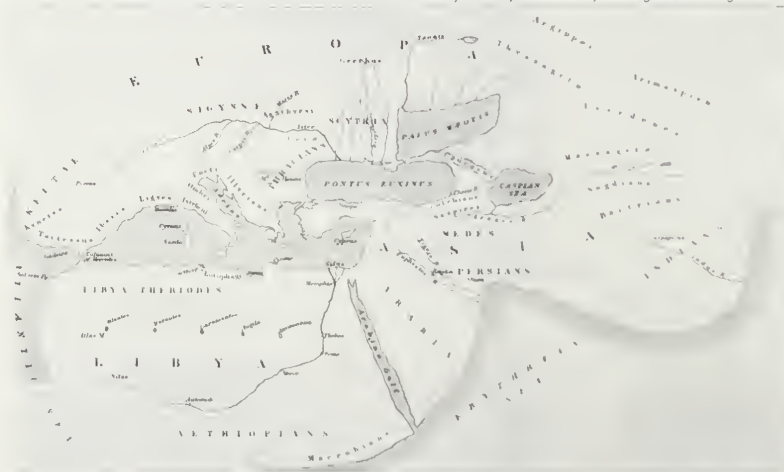


Figure 1: Herodotus' map of the world.

Maps in the New World

fundamental principles of cartography. Ptolemy's map of the world as it was then known marked the culmination of Greek cartography as well as a compendium of accumulated knowledge of the Earth's features at that time (Figure 2).

The Roman period. Although Ptolemy lived and worked at the time of Rome's greatest influence, he was a Greek and essentially a product of that civilization, as was the great library at Alexandria. His works greatly influenced the development of geography, which he defined in map-making terms: "representation in picture of the whole known world, together with the phenomena contained therein." This had considerable influence in directing scholars toward the specifics of map construction and away from the more abstract and philosophical aspects of geography.

One fundamental error that had far-reaching effects was attributed to Ptolemy—an underestimation of the size of the Earth. He showed Europe and Asia as extending over half the globe, instead of the 130 degrees of their true extent. Similarly, the span of the Mediterranean ultimately was proved to be 20 degrees less than Ptolemy's estimate. So lasting was Ptolemy's influence that 13 centuries later Christopher Columbus underestimated the distances to Cathay and India partly from a recapitulation of this basic error.

A fundamental difference between the Greek and Roman philosophies was indicated by their maps. The Romans were less interested in mathematical geography and tended toward more practical needs for military campaigns and provincial administration. They reverted to the older concepts of a disk-shaped world for maps of great areas because they met their needs and were easier to read and understand.

The Roman general Marcus Vipsanius Agrippa, prior to Ptolemy's time, constructed a map of the world based on surveys of the then-extensive system of Roman military roads. References to many other Roman maps have been found, but very few actual specimens survived the Dark

Ages. It is quite probable that the Peutinger Table, a parchment scroll showing the roads of the Roman world, was originally based on Agrippa's map and subjected to several revisions through medieval times.

The tragic turn of world events during the first few centuries of the Christian Era wrought havoc to the accumulated knowledge and progress of mankind. As with other fields of science and technology, progress in geography and cartography was abruptly curtailed. After Ptolemy's day there even appears to have been a retrogression, as exemplified by the Roman trend away from the mathematical approach to mapping.

Great accumulations of documents and maps were destroyed or lost, and the survival of a large part of Ptolemy's work was probably due to its great prestige and popularity. The only other major work on mapping to survive was Strabo's earlier treatise, albeit with some changes from re-copying. Few of the maps and related works of the ancient world have come down to us in their original forms. The tendencies to revise and even recapitulate, when copying manuscripts, are readily understood. Doubtless, the factual content was improved more often than not, but a residual confusion remains when the specimen at hand may be either a true copy of an ancient document or a medieval scholar's version of the subject matter.

THE MIDDLE AGES

Progress in cartography during the early Middle Ages was slight. The medieval mapmaker seems to have been dominated by the church, reflecting in his work the ecclesiastical dogmas and interpretations of Scripture. In fact, during the 6th century Constantine of Antioch created a "Christian topography" depicting the Earth as a flat disk. Thus the Roman map of the world, along with other concepts, continued as authoritative for many centuries. A contemporary Chinese map shows that country occupying most of the world, while the Roman Empire dominates most other maps produced during early Christian times.

Later medieval mapmakers were clearly aware of the

By courtesy of the Library of Congress, Washington, D.C

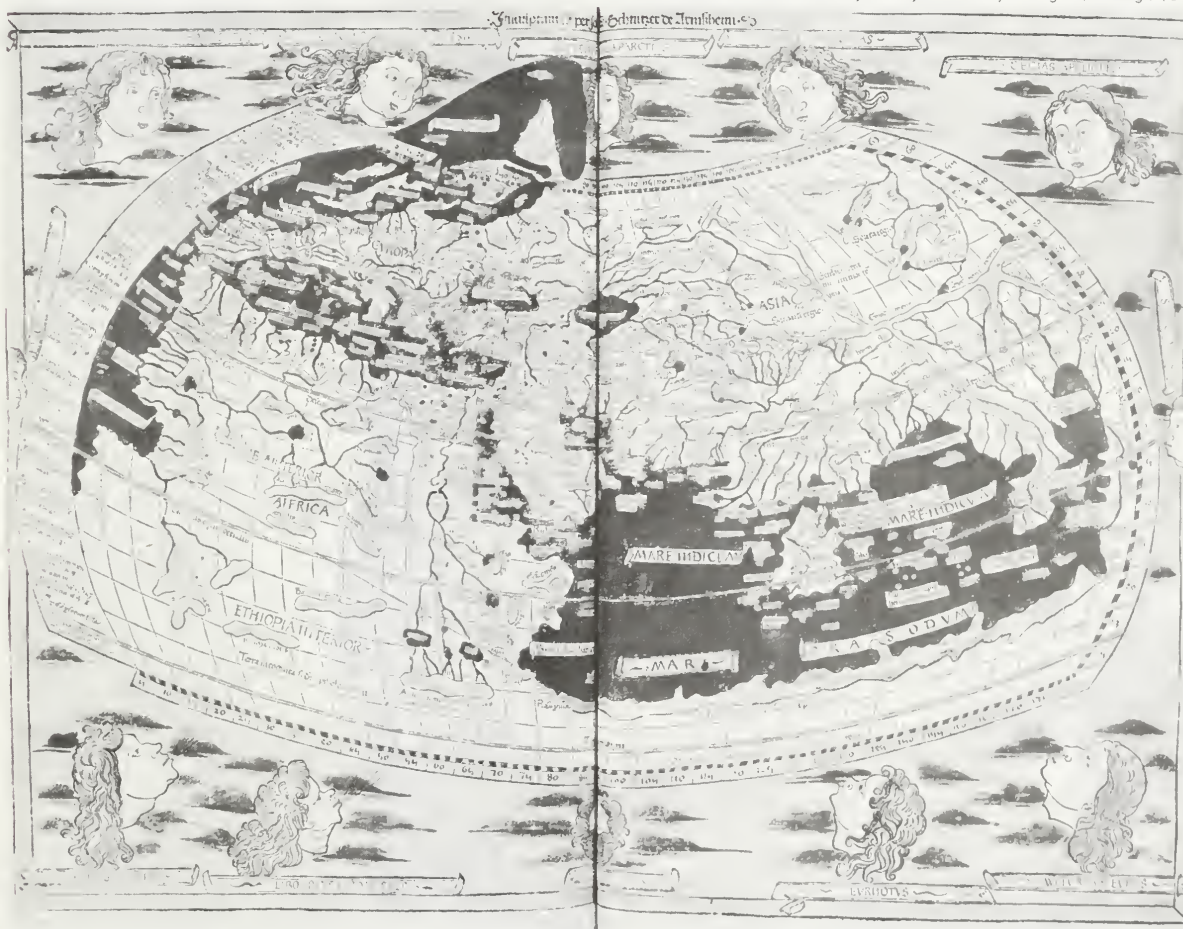


Figure 2: Ptolemy's map of the world, as printed at Ulm, 1482.

Ptolemy's
fundamental
error

Earth's sphericity, but for the most part, maps remained small and schematic, as exemplified by the T and O renderings, so named from the stylized T-form of the major water bodies separating the continents and the O as the circumfluent ocean surrounding the world. The orientation with east at the top of the map was often used, as the word (orientation) suggests.

The earliest navigators coasted from headland to headland; they did not require charts until adoption of the magnetic compass made it possible to proceed directly from one port to another. The earliest record of the magnetic compass in Europe (1187) is followed within a century by the earliest record of a sea chart. This was shown to Louis IX, king of France, on the occasion of his participation in the Eighth Crusade in 1270. The earliest surviving chart dates from within a few years of this event. Found in Pisa and known as the Carta Pisana, it is now in the Bibliothèque Nationale, Paris. Thought to have been made about 1275, it is hand drawn on a sheepskin and depicts the entire Mediterranean Sea. Such charts, often known as portolans named for the portolano or pilot book, listing sailing courses, ports, and anchorages, were much in demand for the increasing trade and shipping. Genoa, Pisa, Venice, Majorca, and Barcelona, among others, cooperated in providing information garnered from their pilots and captains. From repeated revisions, and new surveys by compass, the portolan charts eventually surpassed all preceding maps in accuracy and reliability. The first portolans were hand drawn and very expensive. They were based entirely on magnetic directions and map projections that assumed a degree of longitude equal to a degree of latitude. The assumption did little harm in the Mediterranean but caused serious distortions in maps of higher latitudes. Development of line engraving and the availability, in the 16th century, of large sheets of smooth-surfaced paper facilitated mass production of charts, which soon replaced the manuscript portolans.

Many specimens of portolan charts have survived. Though primarily of areas of the Mediterranean and Black Sea, some covered the Atlantic as far as Ireland, and others the western coast of Africa. Their most striking feature is the system of compass roses (Figure 3), showing directions from various points, and lines showing shortest navigational routes.

Another phenomenon of the late Middle Ages was the great enthusiasm generated by the travels of Marco Polo in the 1270s and 1280s. New information about faraway places, and the stimulation of interest in world maps,

promoted their sale and circulation. Marco Polo's experiences also kindled the desire for travel and exploration in others and were, perhaps, a harbinger of the great age of discovery and exploration.

During Europe's Dark Ages Islamic and Chinese cartography made progress. The Arabs translated Ptolemy's treatises and carried on his tradition. Two Islamic scholars deserve special note. Ibn Haukal wrote a *Book of Ways and Provinces* illustrated with maps, and al-Idrisi constructed a world map in 1154 for the Christian king Roger of Sicily, showing better information on Asian areas than had been available theretofore. In Baghdad astronomers used the compass long before Europeans, studied the obliquity of the ecliptic, and measured a part of the Earth's meridian. Their sexagesimal (based on 60) system has dominated cartography since, in the concept of a 360-degree circle.

Mapmaking, like so many other aspects of art and science, developed independently in China. The oldest known Chinese map is dated about 1137. Most of the area that is now included in China had been mapped in crude form before the arrival of the Europeans. The Jesuit missionaries of the 16th century found enough information to prepare an atlas, and Chinese maps thereafter were influenced by the West.

THE AGE OF DISCOVERY AND EXPLORATION

Revival of Ptolemy. The fall of Byzantium sent many refugees to Italy, among them scholars who had preserved some of the old Greek manuscripts, including Ptolemy's *Geography*, from destruction. The rediscovery of this great work came at a fortunate time because the recent development of a printing industry capable of handling map reproduction made possible its circulation far beyond the few scholars who otherwise would have enjoyed access to it. This, together with a general reawakening of scholarship and interest in exploration, created a golden era of cartography.

The *Geography* was translated into Latin about 1405. Although it had not been completely lost (the Arabs had preserved portions of it), recovery of the complete work, with maps, greatly stimulated general interest in cartography. About 500 copies of the *Geography* were printed at Bologna in 1477, followed by other editions printed in Germany and Italy. The printing process, in addition to permitting the wide diffusion of geographic knowledge, retained the fidelity of the original works. By 1600, 31 Latin or Italian editions had been printed.

Maps of the discoveries. Progress in other technologies

Archivio di Stato di Firenze. Carte nautiche n. 1. Autorizzazione del Ministero dell'Interno, Direzione Generale degli Archivi di Stato (Italia), parere n. 526 del 1971

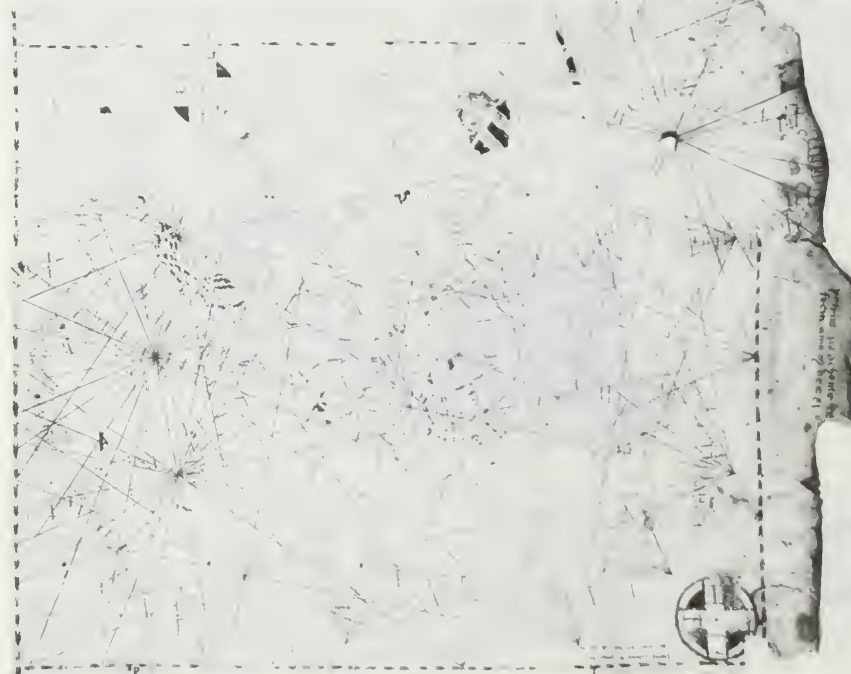


Figure 3: Portolan chart of the eastern Mediterranean area by Petrus Vesconte, 1311.

Portolans

Printing of Ptolemy's *Geography*

such as navigation, ship design and construction, instruments for observation and astronomy, and general use of the compass tended continuously to improve existing map information, as well as to encourage further exploration and discovery. Accordingly, geographic knowledge was profoundly increased during the 15th and 16th centuries. The great discoveries of Columbus, da Gama, Vespucci, Cabot, Magellan, and others gradually transformed the world maps of those days. "Modern" maps were added to later editions of Ptolemy. The earliest was a map of northern Europe drawn at Rome in 1427 by Claudius Claussön Swart, a Danish geographer. Cardinal Nicholas Krebs drew the first modern map of Germany, engraved in 1491. Martin Waldseemüller of St. Dié prepared an edition with more than 20 modern maps in 1513. Maps showing new discoveries and information were at last transcending the classical treatises of Ptolemy.

The most important aspect of postmedieval maps was their increasing accuracy, made possible by continuing exploration. Another significant characteristic was a trend toward artistic and colourful rendition, for the maps still had many open areas in which the artist could indulge his imagination. The cartouche, or title block, became more and more elaborate, amounting to a small work of art. Many of the map editions of this age have become collector's items. The first map printings were made from woodcuts. Later they were engraved on copper, a process that made it possible to reproduce much finer lines. The finished plates were inked and wiped, leaving ink in the cut lines. Dampened paper was then pressed on the plate and into the engraved line work, resulting in very fine impressions. The process remained the basis of fine map reproduction until the comparatively recent advent of photolithography.

The *Cosmographiae*, textbooks of geography, astronomy, history, and natural sciences, all illustrated with maps and figures, first appeared in the 16th century. One of the earliest and best known was that of Petrus Apianus in 1524, the popularity of which extended to 15 more editions. That of Sebastian Münster, published in 1544, was larger and remained authoritative and in demand until the end of the century, reflecting the general eagerness of the times for learning, especially geography.

The foremost cartographer of the age of discovery was Gerhard Kremer, known as Gerardus Mercator, of Flanders. Well educated and a student of Gemma Frisius of Louvain, a noted cosmographer, he became a maker of globes and maps. His map of Europe, published in 1554, and his development of the projection that bears his name made him famous. The Mercator projection solved an age-old problem of navigators, enabling them to plot bearings as straight lines.

Other well-known and productive cartographers of the Dutch-Flemish school are Abraham Ortelius of Antwerp, who prepared the first modern world atlas in 1570, and Jodocus Hondius. Early Dutch maps were among the best for artistic expression, composition, and rendering. Juan de la Cosa, the owner of Columbus' flagship, *Santa María*, in 1500 produced a map recording Columbus' discoveries, the landfall of Cabral in Brazil, Cabot's voyage to Canada, and da Gama's route to India. The first map showing North and South America clearly separated from Asia was produced in 1507 by Martin Waldseemüller. An immense map, 4½ by 8 feet (1.4 by 2.4 metres), printed in 12 sheets, it is probably the first map on which the name America appeared, indicating that Waldseemüller was impressed by the account written by the Florentine navigator Amerigo Vespucci.

In 1529 Diego Ribero, cosmographer to the king of Spain, made a new chart of the world on which the vast extent of the Pacific was first shown. Survivors of Magellan's circumnavigation of the world had arrived in Seville in 1522, giving Ribero much new information.

The first known terrestrial globe that has survived was made by Martin Behaim at Nürnberg in 1492. Many others were made throughout the 16th century. The principal centres of cartographic activity were Spain, Portugal, Italy, the Rhineland, the Netherlands, and Switzerland. England and France, with their growing maritime and

colonial power, were soon to become primary map and chart centres. Capt. John Smith's maps of Virginia and New England, the first to come from the English colonies, were published in London in 1612.

18TH CENTURY TO THE PRESENT

A reformation of cartography that evolved during the 18th century was characterized by scientific trends and more accurate detail. Monsters, lions, and swash lines disappeared and were replaced by more factual content. Soon the only decorative features were in the cartouche and around the borders. The map interiors contained all the increasing information available, often with explanatory notes and attempts to show the respective reliabilities of some portions.

Where mapmakers formerly had sought quick, profitable output based on information obtained from other maps and reports of travelers and explorers, the new French cartographers were scientists, often men of rank and independent means. For expensive ventures, such as the triangulation of two degrees of a meridian to determine the Earth's size more accurately, they were subsidized by the king or the French Academy. Similar trends were developing across Europe.

The new cartography was also based on better instruments, the telescope playing an important part in raising the quality of astronomical observations. Surveys of much higher accuracy were now feasible. The development of the chronometer (an accurate timepiece) made the computation of longitude much less laborious than before; much more information on islands and coastal features came to the map and chart makers.

The rise of national surveys. The development in Europe of power-conscious national states, with standing armies, professional officers, and engineers, stimulated an outburst of topographic activity in the 18th century, reinforced to some extent by increasing civil needs for basic data. Many countries of Europe began to undertake the systematic topographic mapping of their territories. Such surveys required facilities and capabilities far beyond the means of private cartographers who had theretofore provided for most map needs. Originally exclusively military, national survey organizations gradually became civilian in character. The Ordnance Survey of Britain, the Institut Géographique National of France, and the Landestopographie of Switzerland are examples.

In other countries, such as the United States, where defense considerations were not paramount, civilian organizations—e.g., the U.S. Geological Survey and the National Ocean Service (originally Survey)—were assigned responsibility for domestic mapping tasks. Only when World War II brought requirements for the mapping of many foreign areas did the U.S. military become involved on a large scale, with the expansion of the Oceanographic Office (Navy), Aeronautical Chart Service (Air Force), and the U.S. Army Topographic command.

Elaborate national surveys were undertaken only in certain countries. The rest of the world remained largely unmapped until World War II. In some instances colonial areas were mapped by military forces, but except for the British Survey of India, such efforts usually provided piecemeal coverage or generalized and sketchy data. Some important national surveys will be outlined briefly.

The work in France was organized by the French Academy, and in 1748 the *Carte géométrique de la France*, comprising 182 sheets, was authorized. Most of the field observations were accomplished by military personnel. The new map of France as a whole, drawn after the new positions had been computed, caused Louis XV to remark that the more accurate data lost more territory than his wars of conquest had gained. Napoleon, an ardent map enthusiast, planned a great survey of Europe on a 1:100,000 scale, which was well under way when he was overthrown.

During the 18th century Great Britain became the foremost maritime power of Europe, and the Admiralty sponsored many developments in charting as well as improvements in navigation facilities. Because of the Admiralty's prestige, other maritime nations accepted its

Application of science to map-making

Surveys of colonies

Mercator

proposal that the prime meridian for longitude reference should pass through Greenwich. Other achievements in early oceanography were Edmond Halley's magnetic chart, which has been continuously revised from new data. Later similar charts for currents, tides, and prevailing winds were developed.

French progress in mapping stimulated the British to undertake a national survey. The Ordnance Survey was organized in 1791, and the first sheet (Kent), on a scale of one inch to the mile, was published in 1801. By mid-century Ireland had been surveyed at six inches to the mile. In 1858 a Royal Commission approved 1-inch, 6-inch, and 25-inch (1:2500) scales for British mapping. An earlier "first" was John Ogilby's *Britannia*, published in 1675, an atlas of road strip maps plotted by odometer and compass, presaging the modern road map.

A survey of Spain was started in the 18th century. Surveys of several German principalities were combined after unification into the Reichskarte at 1:100,000 scale. A topographic survey of Switzerland was begun in 1832. An Austrian series was started in 1806, from which the Spezialkarte, later considered the most detailed maps of Europe, were derived. In China, under the Communist regime, survey and cartography groups have provided coverage of much of the country with a new 1:50,000-scale map series. Japan established an Imperial Land Survey in 1888, and by 1925 topographic coverage of the home islands, at a scale of 1:50,000, was complete.

International Map of the World (IMW). The International Geographical Congress in 1891 proposed that the participating countries collaborate in the production of a 1:1,000,000-scale map of the world. Specifications and format were soon established, but production was slow in the earlier years since it was first necessary to complete basic surveys for the required data, and during and after World War II there was little interest in pursuing the project. The intention to complete the series was reestablished, however, and many countries have returned to the task. By the mid-1980s the project was nearing completion.

World War II and after. World War I, and to a much greater extent World War II, brought great progress in mapping, particularly of the unmapped parts of the Earth; an appraisal by the U.S. Air Force indicated that in 1940 less than 10 percent of the world was mapped in sufficient detail for even the meagre requirements of pilot charts. A major program of aerial photography and reconnaissance mapping, employing what became known as the trimetrogon method, was developed. Vast areas of the unmapped parts of the world were covered during the war years, and the resulting World Aeronautical Charts have provided generalized information for other purposes since that time. Many countries have used the basic data to publish temporary map coverage until their more detailed surveys can be completed.

The Cold War atmosphere of the 1940s and '50s promoted a continuation of militarily oriented mapping. Both NATO and Warsaw Pact countries continued to improve their maps; NATO developed common symbols, scales, and formats so that maps could be readily exchangeable between the forces of member countries. Postwar economic development programs, in which maps were needed for planning road, railroad, and reservoir constructions, also stimulated much work. The United Nations provides advisory assistance in mapping to countries wishing it.

Among other collaborations, the Inter-American Geodetic Survey, in which the U.S. Army provides instruction and logistic support for mapping, was organized. Although this cooperation primarily involved Latin-American countries, similar arrangements were made with individual countries in other parts of the world. Cooperation and exchange of data in hydrographic surveys, aeronautical charting, and other fields has continued.

Although some terrain data are available for practically all of the world, the data for many sectors remain sketchy. Surveys of Antarctica by the several countries active there are in progress, but the continent will not be completely mapped for some years. The goal of most countries is to achieve adequate coverage for general development needs. Much remains to be done. Even in countries like

the United States that have not yet completed the initial coverage, many of the maps prepared in earlier years are already in need of revision. Thus, even when mapping is completed, requirements for greater detail and revision will continue to make demands upon the funds available.

Aerial photography, which permits accurate and detailed work within feasible cost ranges, has dominated basic mapping in recent years. During World War I aerial photography was used for reconnaissance mapping, and after the war rapid progress was made in optics, cameras, plotting devices, and related equipment. By World War II much of the highly sophisticated equipment now in use had been designed. Electronic distance-measuring devices have made field surveys easier and more accurate, while much improved circle graduation has made theodolites (transits) lighter as well as more precise. Computers and automation, which together have transformed the mapping procedures of yesterday, are described below in the section *Modern mapmaking techniques*.

Mapmaking

ELEMENTS

Map design is a twofold process: (1) the determination of user requirements, with attendant decisions as to map content and detail, and (2) the arrangement of content, involving publication scale, standards of treatment, symbolizations, colours, style, and other factors. To some extent user requirements obviously affect standards of treatment, such as publication scale. Otherwise, the latter elements are largely determined on the basis of efficiency, legibility, aesthetic considerations, and traditional practices.

In earlier productions by individual cartographers or small groups, personal judgments determined the nature of the end product, usually with due respect for conventional standards. Map design for large programs, such as the various national map series of today, is quite formal by comparison. In most countries, the requirements of official as well as private users are carefully studied, in conjunction with costs and related factors, when considering possible changes or additions to the current standards.

Requirements of military agencies often have a decisive influence on map design, since it is desirable to avoid the expense of maintaining both civil and military editions of maps. International organizations and committees are additional factors in determining map design. The fact that development of changes in design and content of national map series may become rather involved induces some reluctance to change, as does the fact that map stocks are usually printed in quantities intended to last for 10 or more years. Also, frequent changes in treatments result in extensive overhauls at reprint time, with consequent inconsistencies among the standing editions.

Planning for the production of a national series involves both technical and program considerations. Technical planning involves the choice of a contour interval (the elevation separating adjacent contour lines, or lines of constant elevation), which in turn determines the height of aerial photography and other technical specifications for each project. The sequence of mapping steps, or operational phases, is determined by the overall technical procedures that have been established to achieve the most efficiency.

The program aspects of planning involve fiscal allotments, priorities, schedules, and related matters.

Production controls also play important roles in large programs, where schedules must be balanced with capacities available in the respective phases to avoid backlogs or dormant periods between the mapping steps. Considering that topographic maps may require two years or more to complete, from authorization to final printing, the importance of careful planning is evident. Many factors, including the weather, can converge to cause delays.

Map scales and classifications. Map scale refers to the size of the representation on the map as compared to the size of the object on the ground. The scale generally used in architectural drawings, for example, is $\frac{1}{4}$ inch to one foot, which means that $\frac{1}{4}$ of an inch on the drawing equals one foot on the building being drawn. The scales of

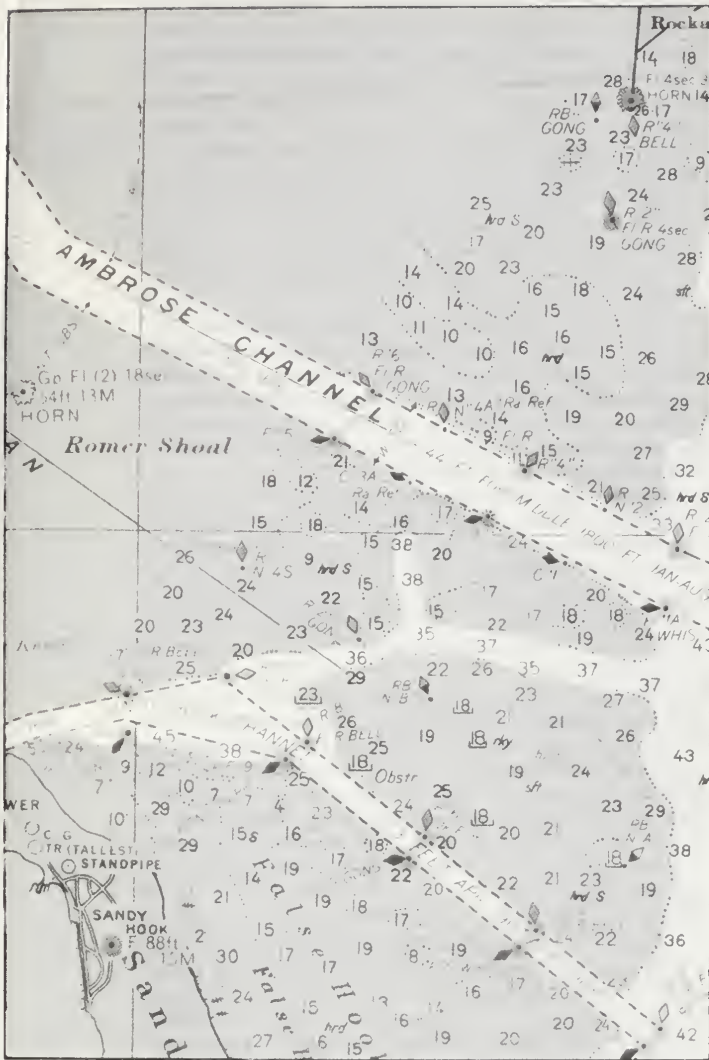


Figure 4: Small section of a large hydrographic chart (U.S. Coast and Geodetic Survey No. 1215) of the approaches to New York harbour. The scale is 1:80,000.

By courtesy of the U.S. Department of Commerce, Environmental Science Services Administration

models of buildings, railroads, and other objects may be one inch to several feet. Maps cover more extensive areas, and it is usually convenient to express the scale by a representative fraction or proportion, as 1/63,360, 1:63,360, or "one-inch-to-one-mile." The scale of a map is smaller than that of another map when its scale denominator is larger: thus, 1:1,000,000 is a smaller scale than 1:100,000. Most maps carry linear, or bar, scales in one or more margins or in the title blocks.

Nautical charts are constructed on widely different scales and can be generally classified as follows: ocean sailing charts are small-scale charts, 1:5,000,000 or smaller, used for planning long voyages or marking the daily progress of a ship. Sailing charts, used for offshore navigation, show a generalized shoreline, only offshore soundings, and are at a scale between 1:600,000 and 1:5,000,000. As an illustration of chart use, a 10-knot ship covers about 29 inches (74 centimetres) at 1:600,000 scale in a day.

General charts, as exemplified by Figure 4, are used for coastwise navigation outside outlying reefs and shoals and are at a scale between 1:100,000 and 1:600,000. Coast charts are intended for use in leaving and entering port or navigating inside outlying reefs or shoals and are at a scale between 1:50,000 and 1:100,000. Harbour charts are for use in harbours and small waterways, with a scale usually larger than 1:50,000.

In rare instances reference may be made to the areal scale of a map, as opposed to the more common linear scale. In such cases the denominator of the fractional reference would be the square of the denominator of the linear scale.

The linear scale may vary within a single map, particu-

larly if the scale is small. Variations in the scale of a map because of the sphericity of the surface it represents may, for practical purposes, be considered as nil. On maps of very large scale, such as 1:24,000, such distortions are negligible (considerably less than variations in the paper from fluctuations of humidity). Precise measurements for engineering purposes are usually restricted to maps of that scale or larger. As maps descend in scale, and distortions inherent to their projection of the spherical surface increase, less accurate measurements of distances may be expected.

Maps may be classified according to scale, content, or derivation. The latter refers to whether a map represents an original survey or has been derived from other maps or source data. Some contain both original and derived elements, usually explained in their footnotes. Producing agencies, technical committees, and international organizations have variously classed maps as large, medium, or small scale. In general, large scale means inch-to-mile and larger, small scale, 1:1,000,000 and smaller, leaving the intermediate field as medium scale. As with most relative terms, these can occasionally lead to confusions but are useful as one practical way to classify maps.

The nature of a map's content, as well as its purpose, provides a primary basis of classification. The terms aeronautical chart, geologic, soil, forest, road, and weather map make obvious their respective contents and purposes. Maps are therefore often classified by the primary purposes they serve. Topographic maps usually form the background for geologic, soil, and similar thematic maps and provide primary elements of the bases upon which many other kinds of maps are compiled.

Map projections. A great variety of map projections has been devised to provide for the various properties that may be desired in maps. In effect, a projection is a systematic method of drawing the Earth's meridians and parallels on a flat surface. Some projections have equal-area properties, while others provide for conformal delineations in which, for small areas, the shape is practically the same as it would be on a globe. Only on a globe can areas and shapes be represented with true fidelity. On flat maps of very large areas, distortions are inevitable. These effects may be minimized by selecting the projection best suited to the purpose of the map to be produced.

Most types of projection can be grouped according to their geometric derivations as cylindrical, conic, or azimuthal (Figure 5). A few cannot be so related or are combinations of these. Terms such as network, graticule, or grid might have been preferable to describe the transposition of meridians and parallels from globe to flat surface, since few systems are actually derived by projection, and most in fact have been formulated by analytic and mathematical processes. The term projection, however, is well established and has some merit in helping the layman to understand the problems and solutions. The theory of trigonometric surveying was disclosed in 1533 by Gemma Frisius, a Flemish mathematician. In 1569 Gerardus Mercator solved the projection problem by producing his famous world map with the meridians vertical and parallels having increased spacing in proportion to the

Types of classification

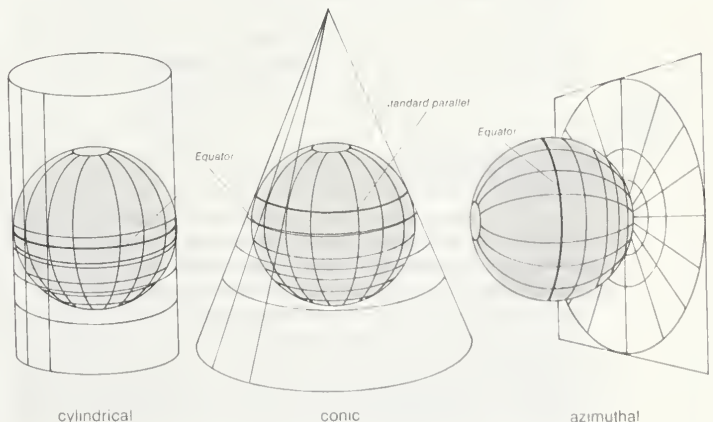


Figure 5: Derivations of basic projections.

secant (a trigonometric function) of the latitude. Edward Wright published mathematical tables (1599) giving the basis of Mercator's projection. Tables for the construction of other commonly used projections have been developed by mapping agencies.

Mercator
projection

Cylindrical projections treat the Earth as a cylinder on which parallels are horizontal lines and meridians appear as vertical lines. The familiar Mercator projection is of this class and has many advantages in spite of the great distortions that it causes in the higher latitudes. Compass bearings may be plotted as straight segments on these projections, which have been traditionally used for nautical charts. On cylindrical projections places of similar latitude appear at the same height. Parallels and meridians may, if desired, be omitted from the body of the map and instead simply indexed at the margins, while lettering can be placed horizontally rather than in a curve. Among the variations of cylindrical projections is the Transverse Mercator, in which the cylinder is tangent to the Earth not along the Equator but along a chosen meridian, a treatment that has advantages in drawing maps that are long in the north-south direction.

Virtually all navigational charts are constructed on the ordinary Mercator projection; the only navigational charts not on ordinary Mercator projections are great-circle charts and charts of the polar regions. Great-circle charts, which are maps of large areas, such as the entire Pacific Ocean, are ordinarily on very small scales with gnomonic projection. The navigator uses them to lay out a track between ports perhaps thousands of miles apart and then transfers the latitudes corresponding, for example, to each 5° of longitude, to his ocean sailing chart. He thus arrives at a series of short rhumb-line courses, each of which makes the same angle with all meridians, that closely approximate the shortest distance between the two ports.

Conic projections are derived from a projection of the globe on a cone drawn with the point above either the North or South Pole and tangent to the Earth at some standard or selected parallel. Occasionally the cone is arranged to intersect the Earth at two closely spaced standard parallels. A polyconic projection, used in large-scale map series, treats each band of maps as part of a cone tangent to the globe at the particular latitude.

Azimuthal, or zenithal, projections picture a portion of the Earth as a flattened disk, tangent to the Earth at a specified point, as viewed from a point at the centre of the Earth, on the opposite side of the Earth's surface, or from a point far out in space. If the perspective is from the centre of the Earth, the projection is called gnomonic; if from the far side of the Earth's surface, it is stereographic; if from space, it is called orthographic (Figure 6).

Modified from I. Fisher and O. M. Miller *World Maps and Globes*

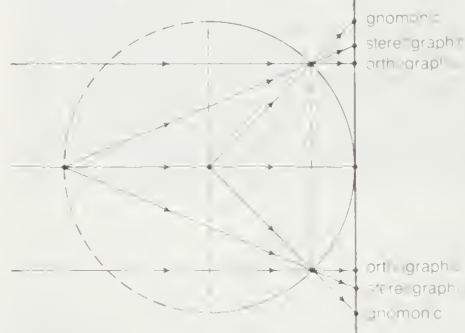


Figure 6: Basis of gnomonic, stereographic, and orthographic projections.

A type of projection often used to show distances and directions from a particular city is the Azimuthal Equidistant. Such measurements are accurate or true only from the selected central point to any other point of interest.

The polar projection is an azimuthal projection drawn to show Arctic and Antarctic areas. It is based on a plane perpendicular to the Earth's axis in contact with the North or South Pole. It is limited to 10 or 15 degrees from the poles. Parallels of latitude are concentric circles, while meridians are radiating straight lines.

Development of reference spheroids. Tables from which map projections of the more familiar kinds may be plotted have been available for some years and have been based on the best determinations of the size and shape of the Earth available at the time of their compilation. The dimensions of Clarke's Spheroid (introduced by the British geodesist Alexander Ross Clarke) of 1866 have been much used in polyconic and other tables. A later determination by Clarke in 1880 reflected the several geodetic surveys that had been conducted during the interim. An International Ellipsoid of Reference was adopted by the Geodetic and Geophysical Union in 1924 for application throughout the world.

Clarke's
Spheroid

The development of electronic distance-measuring systems has facilitated geodetic surveys. During the late 20th century, satellite observation and international collaborations have led to an accurate determination of the size and shape of the Earth and to the possibility of adjusting all existing primary geodetic surveys and astronomical observations to a single world datum.

Geographic and plane coordinate systems. The standard geographic coordinate system of the world involves latitudes north or south of the Equator and longitudes east or west of the Prime Reference Meridian of Greenwich. Map and control point references are stated in degrees, minutes, and seconds carried to the number of decimal places commensurate with the accuracy to which locations have been established.

Geodetic surveys, being of extensive areas, must be adjusted for the Earth's curvature, and reductions must be made to mean sea level for scale. The computations are therefore somewhat involved. As a convenience for engineers and surveyors, many countries have established official plane coordinate systems for each province, state, or sector thereof. By this means, all surveys can be "tied" to control points in the system without transposition to geographic coordinates.

In large countries such as the United States, two basic projections are commonly selected to provide systems with minimum distortions for each state or region. For those long in north-south dimension, the Transverse Mercator is generally used, while for those long in east-west direction, the Lambert conformal (intersecting cone) projection is usually employed. In the case of large regions, two or more zones may be established to limit distortions. Positions of geodetic control points have been computed on the plane coordinate systems and have been made available in published lists.

Basic data for compilation. Maps may be compiled from other maps, usually of larger scale, or may be produced from original surveys and photogrammetric compilations. The former are sometimes referred to as derived maps and may include information from various sources, in addition to the maps from which they are principally drawn. Most small-scale series, such as the International Map of the World and World Aeronautical Charts, are compiled from existing information, though new data are occasionally produced to strengthen areas for which little or doubtful information exists. Thus compiled maps may contain fragments of original information while those representing original surveys may include some existing data of higher order, such as details from a city plat.

Road maps, produced by the millions, are compiled from road surveys, topographic maps, and aerial photography. City maps often represent original surveys, made principally to control engineering plans and construction. Some are, however, compiled from enlargements of topographic maps of the area.

Road and
city map
compila-
tions

Notations regarding the sources from which they were drawn are usually carried on compiled maps. This sometimes includes a reliability diagram showing the areas for which good information was available and those that may be less dependable. Comments regarding certain features or areas, which the editor may deem helpful to the user, may be made in the map itself.

Maps reflecting original surveys, such as a national topographic map series, carry standard marginal information. Date of aerial photography, process and instrumentation employed, notes regarding control and projection, date of

field edit, and other information may be included. References to the availability of adjoining maps and those of other scales or series may also be included. Marginal ticks for intervals of plane coordinate systems, military grids, and other reference features are also shown and appropriately labeled.

Symbolization. Symbols are the graphic language of maps and charts that has evolved through generations of cartographers. The symbols doubtless had their origins as simple pictograms that gradually developed into the conventions now generally used.

Early cartographers recognized that common usages and conventions would minimize confusion and to some extent simplify compilation and engraving. Efforts in this direction were made over the years, but cartographers, being artists of a sort, preferred to vary their styles, and effective standardization was not achieved until comparatively recent times. National agencies in most countries established conventions with due regard to practices in other countries. International Map of the World agreements, NATO conventions, and the efforts of the United Nations and of international technical societies aid standardization.

Symbols may be broadly classed as planimetric or hypsographic or may be grouped according to the colours in which they are conventionally printed. Black is used for names and culture, or works of man; blue for water features, or hydrography; brown for relief, or hypsography; green for vegetation classifications; and red for road classes and special information. There are variations, however, particularly in special-purpose series, such as soil and geologic maps. Symbols will also vary, perforce, because of limitations of space in the smaller scales and the feasibility of drawing some features to true scale on large maps. Legends explain the less obvious symbols on many maps, while explanatory sheets or booklets are available for most standard series, providing general data as well as symbol information. When less familiar symbols are used on maps they are often labeled to prevent misunderstanding. The general located-object symbol, with label, is often used in preference to specific symbols for such objects as windmills and lookout towers for similar reasons.

Planimetric features (those shown in "plan," such as streams, shorelines, and roads) are easier to portray than shapes of land and heights above sea level. Mountains were shown on early maps by sketchy lines simulating profile or perspective appearance as envisioned by the cartographer. Little effort was made at true depiction as this was beyond the scope of available information and existing capabilities. Form lines and hachures, among other devices, were also used in attempting to show the land's shape. Hachures are short lines laid down in a pattern to indicate direction of slope. When it became feasible to map rough terrain in more detail, hachuring developed into an artistic speciality. Some hachured maps are remarkable for their detail and fidelity, but much of their quality depends on the skill of draftsman or engraver. They are little used now, except where relief is incidental.

Contours are by far the most common and satisfactory means of showing relief (see Figure 7). Contours are lines that connect points of equal elevation. The shorelines of lakes and of the sea are contours. Such lines were little used until the mid-19th century, mainly because surveys had not generally been made in sufficient detail for them to be employed successfully. Mean sea level is the datum to which elevations and contour intervals are generally referred. If mean sea level were to rise 20 feet (six metres) the new shoreline would be where the 20-foot contour line is now shown (assuming that all maps on which it is delineated are reasonably accurate).

The quality of contour maps, until recent times, depended largely on the sketching skill of the topographer. In earlier days funds available for topographic mapping were limited, and not much time could be spent in accurate placement. Later, the accurate location of more control points became feasible. An approximate scale of reliability is therefore indicated by the date of a topographic survey, taking into account the respective situations that existed in various countries. Modern surveys, being based on aerial photos and accurate plotting instruments, are generally



Figure 7: Landscape shown in perspective A is represented with the aid of contour lines on map B.

By courtesy of U S Geological Survey

better in detail and accuracy than earlier surveys. The personal skill of individual topographers, long a factor in map evaluations, has therefore been substantially eliminated.

Hill shading, or shaded relief, layer or altitude tinting, and special manipulations of contouring are other methods of indicating relief. Hill shading requires considerable artistry, as well as the ability to visualize shapes and interpret contours. For a satisfactory result, background contours are a necessary guide to the artist. Hypsographic tinting is relatively easy, particularly since photomechanical etching and other steps can be used to provide negatives for the respective elevation layers. Difficulty in the reproduction process is sometimes a deterrent to the use of treatments involving the manipulation of contours.

In the past, three-dimensional maps were laboriously constructed for studies in military tactics and for many other purposes. They were costly to produce, as contour layers had to be cut and assembled, filled with plaster and painted, after which streams, roads, etc., had to be drawn on the surface. Lettering then was applied, and models of large structures, such as buildings and bridges, were added. In view of the time and cost involved in such productions, they were sparingly used until recent years when better production methods and materials became available. During and after World War I a process was developed and improved whereby an aluminum sheet was "raised" by tapping along the contours copied on its surface. When the contours selected for tapping were completed, the sheet became, in effect, a mold for shaping plastic sheets to its convolutions. The map was printed on plastic sheets prior to the thermal process of shaping them to the mold. Sets of relief maps were soon produced in this manner for use in schools, military briefings, and many other activities.

During and after World War II the production of plastic relief maps was greatly expanded, while the processes and equipment were further improved and refined. Most significant among these developments was a pantograph-router, which cuts a model from plaster or other suitable material as the selected contours are followed by the operator on a topographic map. This eliminated the distortions inherent in shaping metal sheets by the tapping process. Selected topographic maps are now published in limited relief editions for military instruction, special displays, and general classroom instruction.

Most relief maps are exaggerated severalfold in the vertical scale. The Earth is remarkably smooth, when viewed in actual scale, and many significant features would hardly be distinguishable on a map without some vertical exaggeration. Mt. Everest, for example, is actually only one-seventh of 1 percent of the Earth's radius in height, or only one-third of an inch (about eight millimetres) at a scale of 1:1,000,000. For this reason relief is usually shown at five, or even 10, times actual scale, depending upon

Hachuring

Pantograph-router

the nature of the area represented. This exaggerated relief scale is always explained in the map legend.

Nomenclature. All possible places and features are identified and labeled to maximize the usefulness of the map. Some names must be omitted, particularly from maps of smaller scales, to avoid overcrowding and poor legibility. The editor must decide which names may be eliminated, while arranging placements so that a maximum number may be accommodated.

Geographic names are the most important, and sometimes the most troublesome, part of the map nomenclature as a whole. Research on existing maps and related documents for a given area may reveal different names for the same features, variations in spelling, or ambiguous applications of names. The field engineer often finds that local usage is confused and sometimes controversial. Various types of official organizations have been established to study the problems submitted and decide the forms and applications that are to be used in government maps and documents. This function is exercised in the United States by the Board on Geographic Names and in the United Kingdom by the Permanent Committee on Geographical Names; worldwide these activities are coordinated by the United Nations Conference on the Standardization of Geographical Names.

Top-
onymics

The science of place-names, or toponymics, has become a significant specialty since World War II, and efforts have been made to establish uniform usages and standards of transliteration throughout the world. Renewed interest in completing the remaining sheets of the International Map of the World, collaborations resulting from military alliances, and efforts of committees of international scientific societies and the United Nations have contributed to these efforts.

At the local levels, however, there are different kinds of problems. The larger scales of most basic topographic map series permit the naming of quite minor hilltops, ridges, streams, and branches, for which designations can be obtained locally. In sparsely settled country few names in actual use may be obtained for minor features, while in other areas inquiries may reveal inconsistencies and confusions in both spelling and application of local names. In some areas, for example, local residents may tend to refer to small streams by the name of the present occupant of the headwater area. The occupants of opposite sides of a mountain sometimes refer to it by different names. In coastal areas the waterman and landsman may use different references for the same features.

A prime opportunity for resolving these problems is presented when a topographic map of an area is prepared for publication. By extensive inquiry and documentation and research of local records and deeds, the appropriate form and application of nearly all names can be determined. Publication and distribution of the map as an official document may then tend to solidify local usage and eliminate the confusions that previously existed.

Lettering

Lettering is selected by the map editor in styles and sizes appropriate to the respective features and the relative importance of each. For topographic maps and most others that follow conventional practice, four basic styles of lettering are used in the Western world. The Roman style is generally used for place-names, political divisions, titles, and related nomenclature. Italic is used for lakes, streams, and other water features. Gothic styles are usually applied to land features such as mountains, ridges, and valleys. Man-made works such as highways, railroads, and canals are usually labeled in slope Gothic capitals, but other distinctive styles are often used for these, together with descriptive notes.

The relative importance of map features is reflected in the different sizes of lettering selected to label them. The most prominent places and features are usually shown in capitals, while lesser ones are labeled with lowercase lettering. In the labeling of cities, however, uppercase lettering is often reserved for state or province capitals. County seats are also labeled in this manner on topographic maps of the United States. For other towns, where lowercase lettering is in order, the sizes selected reflect their relative importance. The use of hand lettering has been aban-

doned in favour of words and figures printed by type or by a photographic process onto transparent material that is "floated" onto the compilation and anchored by an adhesive wax backing in the proper place. Compass roses and graphic scales are added in the same manner.

TYPES AND USES OF MAPS AND CHARTS

World status of mapping and basic data. Before World War I only a few countries, such as Great Britain, France, and Germany, had detailed maps covering their whole national areas. Now many countries have completed coverage of their territories, while others have carried out small-scale coverage and are beginning engineering surveys in selected areas.

It has been demonstrated that the full potential of map usage in a country, state, or province is not realized until some time after complete coverage has become available. When a modern, detailed map replaces an earlier issue, annual distribution can increase dramatically.

Topographic maps provide the basic data for many other kinds as well as working bases for thematic maps showing geology, soils, and vegetation types. The progress of such mapping in the various parts of the world is therefore a primary indicator of the status of cartography in general. A United Nations survey of the status of world mapping is taken periodically. Inquiries are made to the mapping organizations of all member countries regarding the extent of their respective map coverage, publication scales, and related data.

About a third of the world's land area is now covered by maps at scales of 1:75,000 and larger. Some of such coverage is culturally obsolescent or of low structural quality. An additional third is covered by medium-scale topographic maps; *i.e.*, up to 1:125,000 (about two miles to the inch). Some of this is inferior coverage at medium scales, lacking in geodetic control and topographic detail. This is the case with much of China, but most of the mapping is quite adequate for purposes of reconnaissance and as source information for smaller scale maps.

This provides a general indication of the relative reliability of data contained in such world series maps as the 1:1,000,000-scale aeronautical charts and International Maps of the World. Areas of doubtful information are left blank or are drawn with broken lines. In spite of this dearth of reliable data, most of the IMW sheets have been compiled, and most of the aeronautical pilotage charts have been published, to provide navigation continuity across water areas as well as over unmapped parts of the world.

In some areas, however, large-scale topographic maps are not required. Australia, for example, has large-scale coverage only of its populated coastal areas in the east; in the Outback areas 1:250,000-scale maps are considered adequate for most needs, and a program for their production is well under way. Likewise, large areas of tundra, as in Siberia, deserts in many parts of the world, and other sparsely populated areas may be adequately served with medium- or small-scale coverage until specific development sites require engineering maps.

Nautical chart coverage of the world leaves much to be desired. Good progress has been made, however, on areas bordering the continents and islands. The Arctic, Antarctic, South Pacific, and South Atlantic oceans are the most deficient in good coverage. The Defense Mapping Agency, through agreement with the British Admiralty and other chart-producing countries, maintains worldwide coverage that is constantly updated. The National Ocean Service (originally Survey) maintains charts of U.S. coastal waters. The International Hydrographic Organization (until 1967 Bureau), based at Monaco, attempts to stimulate cooperation in improvement of hydrographic data in general. This organization's General Bathymetric Chart of the Oceans shows existing knowledge and is revised from time to time as new data are accumulated.

Coverage of reliable aeronautical charts parallels the availability of topographic maps that provide the essential terrain and cultural data. For this purpose, good 1:250,000-scale maps contain sufficient information for clearance safety and position identification.

Inter-
national
Hydro-
graphic
Organ-
ization

Until recently the progress of geodetic triangulation, the basic survey method, was more or less limited to areas either covered by good topographic maps or scheduled for mapping. Preparations for cadastral surveys, where land partition problems abound, have occasionally led to early geodetic programs. Coastal and other surveys also require good basic control to be fully effective; however, it is again the developed and heavily populated areas that are encompassed with the best geodetic surveys. Electronic distance-measuring systems accelerated the progress of geodetic surveys during the 1960s and extended continental schemes over many ocean areas. International cooperation on satellite triangulation is now in progress, with the prospect that existing triangulation of the continents may soon be tied together and adjusted into a single world datum. The Inter-American Geodetic Survey has made progress in the Americas.

In addition to other applications, aerial photographs provide a useful supplement to topographic maps. Indeed, where maps are not available, aerial photographs invariably serve as map substitutes in spite of inherent distortions and lack of elevation data. Most of the world is covered by aerial photography.

During World War II the U.S. Air Force photographed vast areas of the world, providing reconnaissance maps that were used as bases for aeronautical charts. Much of this information now forms the basis for small-scale map coverage in still remote areas. The system of photography and mapping became known as the trimetrogon process. In it, three wide-angle cameras are used to photograph the terrain from horizon to horizon across the line of flight from an elevation of 20,000 feet (6,100 metres). Detail is usually discernible and plottable for several miles on each side of the line of flight, and occasional points, required for photo-triangulation, can be identified farther out. With higher flight capabilities, wider-angle cameras, and lenses of fine resolution, the progress of aerial photography has been accelerated. Films have been much improved for fineness of emulsion grain and scale stability. Satellite photography and high-altitude flights with super-wide-angle cameras are now under way in the remaining areas of the world. Infrared and colour film developments have greatly improved photo-interpretation capabilities, providing much better delineations for coastal charts, geologic maps, timber and soil classifications, and other thematic mapping.

Types of maps and charts available. Although the range of maps and charts now available in many countries is so extensive that a complete listing is impractical, any list of the principal types would have to include aeronautical (worldwide and national), congressional or political districts, population distribution, geologic (various scales), highways (national and secondary political units), historical, hydrographic (coastal areas, inland waters, foreign waters), national forests, forest types, public land survey plats, soil, and topographic (national and foreign).

The National Atlas of the United States of America, published by the Geological Survey in 1970, contains contributions from all of that country's mapping agencies. Summaries are provided of all thematic and economic data of interest. The atlas also indicates where more detailed information or large-scale specialized maps may be obtained. Many countries have centres where detailed information on existing map series and related data may be

obtained. In the United States this service is performed by the Map Information Office of the U.S. Geological Survey, which publishes and distributes indexes of each state showing map coverage and ordering information. Summary data on geodetic control and aerial photography are also maintained.

The situation is less complex in other countries where mapping activities are concentrated in one or two organizations—e.g., Ordnance Survey in Great Britain and Institut Géographique National in France. The main agencies can advise where maps produced by others may be obtained. Technical societies maintain large map reference libraries and are prime sources of information, as are the map sections of national libraries and museums.

Government and other mapping agencies. The following are the primary agencies of selected countries having advanced mapping programs.

Military agencies play large roles in the mapping activities of many countries. Frequently, a small cadre of officers administers the mapping facilities, while most of the production personnel are civilian. Many countries, such as Iran and the United States, have both civilian and military organizations that collaborate in developing their respective programs and in performing the actual mapping.

Most countries have private and commercial organizations that produce maps. The widely distributed road maps noted earlier are printed by a few large producers who, in cooperation with others, compile the maps. Very large-scale maps—for example, for road construction and other engineering works—are produced under contract by a number of mapping companies. Some local highway departments have their own photogrammetric units to provide or supplement such productions. City surveys and maps for real-estate developments, tax records, power lines, and so on are largely produced by commercial organizations.

Large societies, such as the American Geographical Society, the National Geographic Society, and the Royal Geographical Society, play important roles in addition to being centres of reference as noted above. The National Geographic Society produces popular small-scale maps of the various regions of the world. The American Geographical Society has compiled many maps, most notably a 1:1,000,000 coverage of Hispanic America on standards similar to those of the International Map of the World. Technical societies, such as the American Congress on Surveying and Mapping, the American Society of Photogrammetry, the American Society of Civil Engineers, and others, lend their support to mapping programs and activities. They issue technical papers and hold frequent meetings where new processes and instrumentation are discussed and displayed. *The Manual of Photogrammetry* and *Journal*, produced by the American Society of Photogrammetry, *Photogrammetria*, published by the International Society for Photogrammetry and Remote Sensing, and *Surveying and Mapping*, published by the American Congress on Surveying and Mapping, are prime examples of important contributions that societies make to the overall progress of mapping.

International organizations. Many societies and other types of organizations are now engaged in activities associated with maps and mapping. In general, they encourage cooperation through meetings and articles in their journals; some are more directly concerned with the dissemination of information on the progress of particular kinds of mapping and charting. Standardizations of map treatments and conventional signs as well as the promotion of progress in technical processes are further objectives of such groups.

The United Nations Office of Cartography plays an important role in all of the activities noted above. It maintains records of progress on the International Map of the World and performs related services formerly handled by the Central Bureau of the IMW. Technical assistance in the development of mapping facilities and programs is provided on request. Occasional regional meetings are arranged for groups of countries having similar problems, while the journal *World Cartography* publishes related papers.

Tri-
metrogon
process

The
work of
geographic
societies

Principal Mapping Agencies of Selected Countries

Australia	Division of National Mapping, Department of National Development
Brazil	Serviço Geográfico do Exército
Canada	Survey, Mapping, and Remote Sensing Sector, Department of Energy, Mines and Resources
Chile	Instituto Geográfico Militar
France	Institut Géographique National
Germany	Institut für Angewandte Geodäsie
Iran	Army Geographic Department and National Geographic Centre
United Kingdom	Ordnance Survey
United States	United States Geological Survey
Russia	Glavnoye Upravleniye Geodezii i Kartografii

The Inter-American Geodetic Survey is a special unit of the U.S. Corps of Engineers organized to forward the completion of geodetic surveys and mapping in the Americas. Through technical training and assistance with programs, geodetic surveys in Central and South America have been greatly advanced in recent years. Training in photogrammetry is offered and has promoted the establishment of mapping facilities and programs in many of the collaborating countries.

The Pan American Institute of Geography and History has sponsored regular meetings and consultations on cartography, much in the manner of scientific societies. The consultations are held in different countries each year.

The International Hydrographic Bureau was founded in 1921 in Monaco, where it has been headquartered through the years. It serves as a clearinghouse for information related to hydrography and charting and maintains a General Bathymetric Chart of the World, which is revised periodically to include data furnished by the maritime nations participating in their programs and conferences. Other organizations that promote progress in the various aspects of mapping and charting are the International Association of Geodesy, the International Cartographic Association, the International Civil Aviation Organization, the International Geographical Union, the International Federation of Surveyors, the International Society for Photogrammetry and Remote Sensing, and the International Union of Geodesy and Geophysics.

MODERN MAPMAKING TECHNIQUES

Compilation from existing materials. The preparation of derived maps—*i.e.*, maps that are compiled from other maps or existing data—involves the search for, and evaluation of, all extant data pertaining to the subject area. Depending on the nature of the map to be compiled, thoroughgoing research includes boundary references, historical records, name derivations, and other materials. Selection of the most authentic items, on the frequent occasions when some ambiguities are detected, requires careful study and references to related materials. The sources finally selected may require some adjustment or compromise in order to fit properly with adjacent data. When it becomes evident that some sources are of questionable reliability, the cartographer explains this in the margin of his compilation. Sometimes this is placed in the body of the map where the doubtful features or delineations are located.

When selected materials have been assembled they are reduced to a common scale and copied on the compilation base, often in differentiating colours for the respective features. Reductions to a common scale are usually made by photography but may be made by projection and traced directly on the drawing. Minor adjustments may have to be made during compilation even though the source materials are of good quality. In particular, the need to make appropriate generalizations, omitting some details in smaller scale maps, requires much study and judgment.

Except for the new methods of preparing final colour-separation plates by scribing (described below), rather than by drafting or copperplate engraving, compilation processes have changed little over the years. Automatic-focusing projectors and better illumination have made the tracing of selected data at compilation scale easier. Better and more extensive facilities for photoreduction and copying, improved light tables, and a wider choice of drafting materials and instruments have served to facilitate compilation. The basic chores of research, selection of best data, and adjustment of these into the compilation, however, remain essentially the same.

The preparation of small-scale maps from large ones is sometimes simpler than the process just described, which pertains to compilation from a miscellany of differing sources. The relatively straightforward preparation of 1:62,500-scale maps from those of the 1:24,000-scale series, for example, may require little more than photoreduction and colour-separation drafting, or scribing. Even in this case some generalizations, as well as omission of a few of the least important details, are in order. To avoid the considerable expense involved in such scale conver-

sions, straight photoreduction of colour-separation plates appears to be a promising procedure.

Larger reductions from one map series to another—1:62,500 to 1:250,000 for instance—are more of a problem, since the need for generalization is greater and the omission of many details is involved. The considerable differences in road and other symbol sizes also create displacement problems.

The component maps are reduced, and the negatives are cut and assembled into a mosaic on a clear sheet of plastic, the master negative of which provides guide copy for the several colour-separation plates required, which are then completed for reproduction. More often, however, it is necessary to make an intermediate compilation rather than burden the draftsman with too many adjustments to be made while following copy on the colour-separation plates. The intermediate scale for initial reduction of the component maps provides better legibility than direct photography to reproduction scale. This negative mosaic is copied on a metal-mounted drafting board. A compiler then inks the whole map, usually in three or more contrasting colours. He also draws roads and other symbols at the intermediate size, so that they will reduce to proper dimensions at reproduction scale, and makes the necessary displacement adjustments. Minor features and terrain details to be omitted on the new map are not inked in. The drawing is now ready for photoreduction to the final colour-separation plates, providing much better copy for the draftsman or engraver than direct reduction in one step would have produced.

Most smaller scale map series are prepared from large-scale maps as described above. In earlier days original reconnaissance surveys were made at small scales such as 1:192,000 for publication at 1:250,000. Ideally, the small-scale series of maps should be compiled progressively from those of larger scale and greater detail. Most countries, however, started their mapping programs with relatively small-scale reconnaissance surveys because of economic considerations. Later, affluence and technical competence permitted mapping at larger scales with better accuracy.

Geologic, soil, and other thematic maps usually have a topographic base from which woodland tints and road classification printings have been omitted. Such a map, therefore, has a topographic background printed in subdued colours on which the geologic or soil patterns are overprinted in prominent colours. Small-scale thematic maps showing weather patterns, vegetation types, and a large amount of economic and other information are of similar origin. Backgrounds are drawn from appropriate outline maps of provinces, countries, or regions of the world, while overlaying subject matters are compiled from specialized sources of information.

The generalization of detail is a problem that frequently confronts the cartographer in original mapping and in reducing the scale of existing maps. There are two principal reasons for taking such liberties (or topographic license in the case of the original mapping). The primary purpose is to avoid overcrowding and the resulting poor legibility. In addition, the degree of generalization or detail should be as consistent as possible throughout the map. Generalizations in some parts and excessive detail in others confuse the user and make the map's reliability suspect. Effective generalization requires good judgment based on seasoned knowledge and experience.

In approaching such problems as the thousands of islets in the Stockholm archipelago or the thousands of small lakes in the Alaskan tundra areas, when the map scale will accommodate only a small number, the cartographer may decide to draw the features in groupings that reflect the patterns shown in the large-scale source maps or aerial photos. This is difficult and at best presents the nature of the respective areas rather than a literal portrayal. There is also the possibility that the source maps may already have been generalized by some omissions to accommodate to their own scales. Another device is to note, in appropriate text or marginal references, that many minor lakes or islets are omitted because of scale. Such areas may also be symbolized and explained. The "pattern" representation noted above is actually a form of symbolization.

General
Bathy-
metric
Chart of
the World

Use of
topo-
graphic
back-
grounds

Straight
photo-
reduction

Intricate coastlines are also extremely difficult to generalize consistently. Here again, the purpose is to omit minor details while retaining the main features and their distinguishing characteristics. These and many equally perplexing questions arise in preparing maps of very small scale from any source. The problems of equalization of detail are also present in such cases. The topographer of earlier days had the equalization problem between areas close at hand and those viewed distantly. In addition, the topographer had to deal with terrain on the far sides of obscuring features.

Photogrammetrists—that is, persons who compile original maps from aerial photos—have similar problems when, for example, one side of a ridge is seen in more detail than the opposite side. Indeed, in steep terrain, parts of the far sides of some mountains are not seen at all. Appropriate steps must be taken in such cases to avoid differing renditions on opposite sides of the mountain. This may be accomplished by adding, in field completion of the manuscript map, the segments not seen by the photogrammetrist; or additional aerial photography, patterned to cover the obscured sectors, may be requested.

Map production from original surveys. The instrumentation, procedures, and standards involved in making original surveys have improved remarkably in recent years. Geodetic, topographic, hydrographic, and cadastral surveys have been facilitated by the application of electronics and computer sciences. At the same time, superior optics and more refined instruments, in general, have enhanced the precision of observations and accuracies of the end products.

The improved quality of surveys has increased the reliability of maps and charts based on them. In turn, the greater output of basic data has accelerated the production of maps and charts, while parallel improvements in processing steps have increased the volume and improved the final product. In a sense the production of maps from original surveys parallels the process steps after a compilation is made from derived sources. This phase is sometimes referred to as map finishing and involves editing, colour separation, and printing. In original surveys for topographic maps and nautical charts, however, the end products are provided for in all the process steps leading to the completed basic manuscript. The manuscript scale is, for example, selected to accommodate the plotting instruments involved as well as the final rendering for printing. In early years it was usual to choose a manuscript scale somewhat larger than that prescribed for publication. This was to allow for some generalization and line refinement in the final reduction. Thus, maps to be published at 1:62,500 scale were plotted in the field at 1:48,000 or thereabouts. With modern photogrammetric instruments, plotting is usually at reproduction scale.

Maps are not directly derived from geodetic surveys, and only land-line plats are produced from cadastral surveys. Accordingly, the primary original map and chart productions are those from topographic and hydrographic surveys. The surveys are somewhat similar as the nautical chart is, in effect, a topographic map of the coast with generalized offshore topography interpolated from depth soundings.

A variety of electronic devices are used to determine a survey ship's precise location while taking soundings, which are also made with electronic equipment. Both hydrographic and topographic surveys now employ aerial photography and precise plotting instruments to develop the base map. In order to simplify the description of modern mapmaking techniques, the process developed for topographic mapping will be described below, with comments where procedures for nautical charts differ significantly. Both processes start by expanding upon the basic control previously established from geodetic surveys.

Surveying, in which the facts are discovered and recorded, must precede mapping, in which the facts are presented in graphic form. Surveying involves (1) global positioning, in which the area to be mapped is located on the Earth's surface, usually by fixing a number of points in the area by astronomical observations or, after the techniques became available, by satellite or radar procedures; (2) establishing

the framework, in which these points, and commonly many others connected by some combination of distance and angle measurements, are integrated into an accurately defined structure—like the steel framework of a modern building—on which the detail survey is based; and (3) making the detail survey, which establishes by less accurate (and therefore cheaper) methods the relative positions and shapes of the features being mapped. Constant reference to the framework prevents the errors in the detail survey from accumulating and growing unacceptably large.

Mapping also consists of three steps: (1) fair drawing, in which the accurate but not publishable surveyor's plot is redrawn by a skilled cartographer with uniform lines and lettering and, if a multicoloured map is being produced, is separated into several drawings, one for each colour; (2) reproduction, in which a negative is prepared from each of the fair-drawn originals and special colouring (to represent areas of vegetation, for example) is added; and (3) printing, in which a printing plate is made from each negative, the plates are mounted on a press, proofs (a few trial copies) are made to facilitate correction of errors and blemishes, and the final maps are produced.

Final steps in map preparation. After all the features visible in the aerial photographs have been mapped, the manuscripts are contact printed on coated plastic sheets for review by the field engineer. He examines the whole map, adding such details as houses, trails, and fences that were not visible or were overlooked by the photogrammetrist. Political lines such as state, county, and township limits are located, as are geographic and other names in local use. Roads are classified, and woodland outlines are checked.

Contour accuracy is tested if the operator has noted areas that may be weak. The determination of names involves extensive local inquiry, as do political lines, and both may require research of records.

In remote areas it is more efficient to combine the above activities with supplemental control survey to avoid the extra field phase. Then the photos must be carefully examined and annotated for the compiler, while buildings must be encircled or pricked. Roads are classified and political lines located in the usual manner and noted on the photos or overlays.

Field corrections are applied to the original manuscripts. They must be scribed (engraved) on the originals so that guide copies can be prepared by contact printing for final colour-separation scribing. At this time all factual detail is carefully checked. Editing may proceed, to conserve time, while the colour-separation scribing is in progress. The editor reviews all names, boundaries, and related data, comparing them to information thereon that may be available from other sources. The editor's function is to see that the map conforms to standard conventions and is clear, legible, and free of errors.

Controversial names, or those found to be in confused or ambiguous spelling or usage, are documented and referred to an appropriate official body. The designation of type styles and sizes as well as placement of lettering is another function of the map editor.

Because modern topographic maps are printed in several colours, separate plates must be prepared for each. Some of the earliest maps were printed from woodcuts, usually in a single colour. Various hand processes were developed through the years, culminating in the fine rendering of copperplate engraving, which dominated the map production industry for many years. The process became obsolete, however, with increased production demands and the development of efficient printing presses. After World War II engraving on glass, and later on coated plastic sheets, was developed to a point that recovered the fineness of copper engraving. These methods of engraving have become firmly established in map production throughout the world.

In the negative engraving or scribing process, guide copy is printed on several sheets of plastic coated with an opaque paint, usually yellow. The scribe follows copy on the respective plates by engraving through the coating. Because arc light can pass only through the engraving scratches, the completed engravings are, in effect, negatives from which

Handling
corrections

Use of
electronic
devices
in ship
surveys

Scribing
process

the press plates are made. The finest lines (0.002 inch, or 0.05 millimetre, wide), such as intermediate contours, are engraved freehand. Heavier lines, such as index contours, engraved at 0.007 inch, may require a small tripod to assure that the scribe is perfectly vertical. Gravers for double-lined roads, others for buildings, and templates, or patterns, for a variety of symbols are used. Woodland and similar boundaries and shorelines are contact printed and etched on their respective coated sheets, and the areas of the woodland or water are then peeled off, leaving open windows for their respective features. If portions of scrub, orchard, or vineyard are contained in the "woodland" plate, negative sections for these are stripped into their respective locations. Press plates are then processed from the negatives.

A combined-colour proof is then made by successively printing the several completed negatives on a sensitized white plastic sheet that serves for the final checking and review of all aspects of the map. After all corrections have been made, the negatives are ready for the reproduction process.

Nearly all maps are now printed by rotary offset presses, using flexible aluminum-alloy printing plates. The system uses surface plates (very slightly raised or recessed) as opposed to the letterpress and intaglio processes, which involve greater image heights and depths respectively. In the printing sequence, ink goes from the plate to a rubber blanket to the paper. Thus, the printing plate is positive, or right-reading, as is the printed map. The negatives from which the printing plates are prepared are accordingly wrong-reading. This is the process for so-called surface plates. To retain fineness of line on very long runs (10,000 or more impressions), some map printers prefer "etched" plates, prepared from film positives. Both may be considered essentially surface plates, however, since the respective raise or recess is quite small.

Presses are of many varieties and makes. Huge multicolour types are used in large plants, printing several colours at a time. In effect, a multicolour press is several presses built into one. Each unit has three cylinders for plate, rubber blanket, and paper as well as rollers for water and ink. Presses with automatic feed may produce as many as 6,000 impressions per hour, while hand-fed types are limited to about 2,500 per hour.

Nautical charts are commonly large, 28 by 40 inches (70 centimetres by 1 metre) being an internationally accepted maximum size. In order that a navigator may work with them efficiently, charts must be kept with a minimum of folding in drawers in a large chart table in a compartment of the ship having ready access to the navigating bridge, known as the chart room or chart house. Such structures are not possible in small craft, which therefore require charts of a more convenient size. With the recognition that there are many more small boats in the world, particularly recreational craft, than there are ships and that they are navigated primarily by piloting rather than by celestial or electronic means, many hydrographic offices have given attention to the production of special chart series in a small format for yacht navigators.

SC charts

A typical series is that produced by the U.S. Coast and Geodetic Survey with the designation SC (for small craft). Such charts are only 15 inches (38 centimetres) in the vertical dimension and thus need to be folded only in the vertical direction. Printed on both sides of the sheet, they are oriented along the most probable route rather than by parallels and meridians. Several are stapled together into a stiff cardboard folder for protection. Along with the ordinary chart information, they contain a year of tide tables and information on small-craft facilities in the area. New editions are produced annually.

Practical uses of charts impose some constraints on the selection of colour. Red, for example, would logically be chosen as the color in which to print warnings of navigational hazards. But navigators, who must work at night, prefer to retain the darkness adaptation of their eyes by viewing their charts under red light. Under such illumination, red, orange, and buff are invisible. Hence these colours have been superseded by magenta, purple, and gray.

Charts are working instruments, and, since ships often voyage far from where replacement charts are readily obtainable, hydrographic offices give attention to the quality of the paper on which charts are printed. A ship's reckoning is kept in pencil and erased after each voyage. Thus, printing stock that permits multiple erasures is chosen. In view of the environment where charts are used, another quality commonly sought is high wet strength.

Automation in mapping. During the past few decades, there was much interest in the automation of mapping processes, and considerable progress was made in this area. Achievements in the fields of electronics, high-speed digital computers, and related technologies provided a favourable period for such progress. In Great Britain, development of a set of procedures utilizing automatic elements, known as the Oxford System, was begun.

Some success was also achieved in the difficult area of automatic plotting. Instruments now available can automatically scan a stereo model and generate approximate profiles from which contours may be interpolated. Some steps, however, must be closely monitored or else performed completely by the operator. Contouring interpolated from a profile scan is inferior to an operator's delineation. This contouring meets some less exacting requirements for elevation data, and refinements in the system are improving its precision. The need for human intervention when automatic devices get "lost" is not a decisive drawback, as one operator can monitor several machines. The reduction of tedious and repetitive steps for stereo-operators offers a significant advancement.

Coordinatographs with high repeat accuracies facilitate the automatic plotting of control points and projection intersections. Line work can also be drafted or scribed automatically by the same process, but the respective features must first be coded to provide the necessary input tapes. Automatic colour scanning and discrimination is operational but has not become widely used; it is still necessary for an operator to trace the various features on the manuscript to code them. Obviously, little is to be gained by automatic scribing until the input can be provided automatically. Coded line work can be displayed on a cathode-ray tube and corrected with a light pen, but it is much simpler to check and correct the manuscript or finished drawing. Systems of automatic type placement at present offer only marginal advantages over conventional methods. In short, automation has made substantial advances but has not become fully operational in a practical sense.

An aspect of automation that is developing rapidly concerns graphic data acquisition, storage, and retrieval. Data banks are being accumulated by specialized users of topographic information, often to produce thematic maps showing soil types, vegetation classifications, and a variety of other information. Such data banks are usually organized in two parts: one for line work, such as boundaries, and the other for descriptive information or classifications. Assuming that the necessary inputs have been made to the data bank, special plats can be generated speedily. Examples of such graphics include profiles showing elevations along a selected radio propagation path and cross sections for earthwork on roads and other construction.

(C.F.F./Ed.)

Surveying

Surveying is a means of relatively large-scale, accurate measurement of Earth surfaces. It includes the determination of the measurement data, the reduction and interpretation of the data to usable form, and, conversely, the establishment of relative position and size according to given measurement requirements. Thus, surveying has two similar but opposite functions: (1) the determination of existing relative horizontal and vertical position, such as that used for the process of mapping, and (2) the establishment of marks to control construction or to indicate land boundaries.

Surveying has been an essential element in the development of the human environment for so many centuries that its importance is often forgotten. It is an imperative

requirement in the planning and execution of nearly every form of construction. Surveying was essential at the dawn of history, and some of the most significant scientific discoveries could never have been implemented were it not for the contribution of surveying. Its principal modern uses are in the fields of transportation, building, apportionment of land, and communications.

Except for minor details of technique and the use of one or two minor hand-held instruments, surveying is much the same throughout the world. The methods are a reflection of the instruments, manufactured chiefly in Switzerland, Austria, Great Britain, the United States, Japan, and Germany. Instruments made in Japan are similar to those made in the West.

HISTORY

It is quite probable that surveying had its origin in ancient Egypt. The Great Pyramid of Khufu at Giza was built c. 2700 BC, 755 feet long and 480 feet high. Its nearly perfect squareness and north-south orientation affirm the ancient Egyptians' command of surveying.

Evidence of some form of boundary surveying as early as 1400 BC has been found in the fertile valleys and plains of the Tigris, Euphrates, and Nile rivers. Clay tablets of the Sumerians show records of land measurement and plans of cities and nearby agricultural areas. Boundary stones marking land plots have been preserved. There is a representation of land measurement on the wall of a tomb at Thebes (1400 BC) showing head and rear chainmen measuring a grainfield with what appears to be a rope with knots or marks at uniform intervals. Other persons are shown. Two are of high estate, according to their clothing, probably a land overseer and an inspector of boundary stones.

There is some evidence that, in addition to a marked cord, wooden rods were used by the Egyptians for distance measurement. They had the groma, which was used to establish right angles. It was made of a horizontal wooden cross pivoted at the middle and supported from above (Figure 8, left). From the end of each of the four arms hung a plumb bob. By sighting along each pair of plumb bob cords in turn, the right angle could be established. The device could be adjusted to a precise right angle by observing the same angle after turning the device approximately 90°. By shifting one of the cords to take up half the error, a perfect right angle would result.

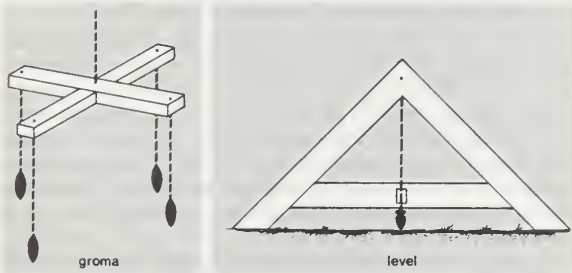


Figure 8: Egyptian surveying instruments.

There is no record of any angle-measuring instruments of that time, but there was a level (Figure 8, right) consisting of a vertical wooden A-frame with a plumb bob supported at the peak of the A so that its cord hung past an indicator, or index, on the horizontal bar. The index could be properly placed by standing the device on two supports at approximately the same elevation, marking the position of the cord, reversing the A, and making a similar mark. Halfway between the two marks would be the correct place for the index. Thus, with their simple devices, the ancient Egyptians were able to measure land areas, replace property corners lost when the Nile covered the markers with silt during floods, and build the huge pyramids to exact dimensions.

The Greeks used a form of log line for recording the distances run from point to point along the coast while making their slow voyages from the Indus to the Persian Gulf about 325 BC. The magnetic compass was brought to the West by Arab traders in the 12th century AD. The astro-

labe was introduced by the Greeks in the 2nd century BC. An instrument for measuring the altitudes of stars, or their angle of elevation above the horizon, took the form of a graduated arc suspended from a hand-held cord. A pivoted pointer that moved over the graduations was pointed at the star. The instrument was not used for nautical surveying for several centuries, remaining a scientific aid only.

During their occupation of Egypt, the Romans acquired Egyptian surveying instruments, which they improved slightly and to which they added the water level and the plane table. About 15 BC the Roman architect and engineer Vitruvius mounted a large wheel of known circumference in a small frame, in much the same fashion as the wheel is mounted on a wheelbarrow; when it was pushed along the ground by hand it automatically dropped a pebble into a container at each revolution, giving a measure of the distance traveled. It was, in effect, the first odometer.

The water level consisted of either a trough or a tube turned upward at the ends and filled with water. At each end there was a sight made of crossed horizontal and vertical slits. When these were lined up just above the water level, the sights determined a level line accurate enough to establish the grades of the Roman aqueducts. In laying out their great road system, the Romans are said to have used the plane table. It consists of a drawing board mounted on a tripod or other stable support and of a straightedge—usually with sights for accurate aim (the alidade) to the objects to be mapped—along which lines are drawn. It was the first device capable of recording or establishing angles. Later adaptations of the plane table had magnetic compasses attached.

Plane tables were in use in Europe in the 16th century, and the principle of graphic triangulation and intersection was practiced by surveyors. In 1615 Willebrord Snell, a Dutch mathematician, measured an arc of meridian by instrumental triangulation.

In 1620 the English mathematician Edmund Gunter developed a surveying chain, which was superseded only by the steel tape in the beginning of the 20th century.

The study of astronomy resulted in the development of angle-reading devices that were based on arcs of large radii, making such instruments too large for field use. With the publication of logarithmic tables in 1620, portable angle-measuring instruments came into use. They were called topographic instruments, or theodolites. They included pivoted arms for sighting and could be used for measuring both horizontal and vertical angles. Magnetic compasses may have been included on some.

The vernier, an auxiliary scale permitting more accurate readings (1631), the micrometer microscope (1638), telescopic sights (1669), and spirit levels (about 1700) were all incorporated in theodolites by about 1720. Stadia hairs were first applied by James Watt in 1771. The development of the circle-dividing engine about 1775, a device for dividing a circle into degrees with great accuracy, brought one of the greatest advances in surveying methods, as it enabled angle measurements to be made with portable instruments far more accurately than had previously been possible.

By the late 18th century modern surveying can be said to have begun. One of the most notable early feats of surveyors was the measurement in the 1790s of the meridian from Barcelona, Spain, to Dunkirk, Fr., by two French engineers, Jean Delambre and Pierre Méchain, to establish the basic unit for the metric system of measurement.

Many improvements and refinements have been incorporated in all the basic surveying instruments. These have resulted in increased accuracy and speed of operations and have opened up possibilities for improved methods in the field. In addition to modification of existing instruments, two revolutionary mapping and surveying changes have been introduced: photogrammetry, or mapping from aerial photographs (about 1920), and electronic distance measurement, including the adoption of the laser for this purpose as well as for alignment (in the 1960s). Important technological developments benefiting surveying in the 1970s included the use of satellites as reference points for geodetic surveys and electronic computers to speed the processing and recording of survey data. (J.Ly./Ed.)

Egyptian and other ancient devices

17th-century portable instruments

The start of modern surveying

Greek log line and Chinese lodestone

MODERN SURVEYING

Basic control surveys. Geodetic surveys involve such extensive areas that allowance must be made for the Earth's curvature. Baseline measurements for classical triangulation (the basic survey method that consists of accurately measuring a base line and computing other locations by angle measurement) are therefore reduced to sea-level length to start computations, and corrections are made for spherical excess in the angular determinations. Geodetic operations are classified into four "orders," according to accuracy, the first-order surveys having the smallest permissible error. Primary triangulation is performed under rigid specifications to assure first-order accuracy.

Efforts are now under way to extend and tie together existing continental networks by satellite triangulation so as to facilitate the adjustment of all major geodetic surveys into a single world datum and determine the size and shape of the Earth spheroid with much greater accuracy than heretofore obtained. At the same time, current national networks will be strengthened, while the remaining amount of work to be done may be somewhat reduced. Satellite triangulation became operational in the United States in 1963 with observations by Rebound A-13, launched that year, and some prior work using the Echo 1 and Echo 2 passive reflecting satellites. The first satellite specifically designed for geodetic work, Pageos 1, was launched in 1966.

A first requirement for topographic mapping of a given area is an adequate pattern of horizontal and vertical control points, and an initial step is the assembly of all such existing information. This consists of descriptions of points for which positions (in terms of latitude and longitude) and elevations above mean sea level have been determined. They are occasionally located at some distance from the immediate project, in which case it is necessary to expand from the existing work. This is usually done on second- or third-order standards, depending upon the length of circuits involved.

The accuracy of survey measurements can be improved almost indefinitely but only at increased cost. Accordingly, control surveys are used; these consist of a comparatively few accurate measurements that cover the area of the project and from which short, less accurate measurements are made to the objects to be located. The simplest form of horizontal control is the traverse, which consists of a series of marked stations connected by measured courses and the measured angles between them. When such a series of distances and angles returns to its point of beginning or begins and ends at stations of superior (more accurate) control, it can be checked and the small errors of measurement adjusted for mathematical consistency. By assuming or measuring a direction of one of the courses and rectangular coordinates of one of the stations, the rectangular coordinates of all the stations can be computed.

A system of triangles usually affords superior horizontal control. All of the angles and at least one side (the base) of the triangulation system are measured. Though several arrangements can be used, one of the best is the quad-

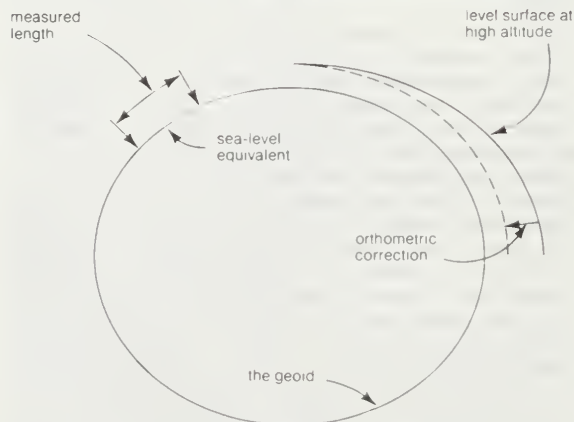


Figure 9: Measured length corrected to its sea-level equivalent (on the geoid) and the orthometric correction for high altitude measurements (see text).

range or a chain of quadrangles. Each quadrangle, with its four sides and two diagonals, provides eight angles that are measured. To be geometrically consistent, the angles must satisfy three so-called angle equations and one side equation. That is to say the three angles of each triangle, which add to 180° , must be of such sizes that computation through any set of adjacent triangles within the quadrangles will give the same values for any side. Ideally, the quadrangles should be parallelograms. If the system is connected with previously determined stations, the new system must fit the established measurements.

When the survey encompasses an area large enough for the Earth's curvature to be a factor, an imaginary mathematical representation of the Earth must be employed as a reference surface. A level surface at mean sea level is considered to represent the Earth's size and shape, and this is called the geoid. Because of gravity anomalies, the geoid is irregular; however, it is very nearly the surface generated by an ellipse rotating on its minor axis—*i.e.*, an ellipsoid slightly flattened at the ends, or oblate. Such a figure is called a spheroid. Several have been computed by various authorities; the one usually used as a reference surface by English-speaking nations is (Alexander Ross) Clarke's Spheroid of 1866. This oblate spheroid has a polar diameter about 27 miles (43 kilometres) less than its diameter at the Equator.

Because the directions of gravity converge toward the geoid, a length of the Earth's surface measured above the geoid must be reduced to its sea-level equivalent—*i.e.*, to that of the geoid (see Figure 9). These lengths are assumed to be the distances, measured on the spheroid, between the extended lines of gravity down to the spheroid from the ends of the measured lengths on the actual surface of the Earth. The positions of the survey stations on the Earth's surface are given in spherical coordinates.

Bench marks, or marked points on the Earth's surface, connected by precise leveling constitute the vertical controls of surveying. The elevations of bench marks are given in terms of their heights above a selected level surface called a datum. In large-level surveys the usual datum is the geoid. The elevation taken as zero for the reference datum is the height of mean sea level determined by a series of observations at various points along the seashore taken continuously for a period of 19 years or more. Because mean sea level is not quite the same as the geoid, probably because of ocean currents, in adjusting the level grid for the United States and Canada all heights determined for mean sea level have been held at zero elevation.

Because the level surfaces, determined by leveling, are distorted slightly in the area toward the Earth's poles (because of the reduction in centrifugal force and the increase in the force of gravity at higher latitudes), the distances between the surfaces and the geoid do not exactly represent the surfaces' heights from the geoid. To correct these distortions, orthometric corrections (see Figure 9) must be applied to long lines of levels at high altitudes that have a north-south trend.

Trigonometric leveling often is necessary where accurate elevations are not available or when the elevations of inaccessible points must be determined. From two points of known position and elevation, the horizontal position of the unknown point is found by triangulation, and the vertical angles from the known points are measured. The differences in elevation from each of the known points to the unknown point can be computed trigonometrically.

The National Ocean Service in recent years has hoped to increase the density of horizontal control to the extent that no location in the United States will be farther than 50 miles (80 kilometres) from a primary point, and advances anticipated in analytic phototriangulation suggest that the envisioned density of control may soon suffice insofar as topographic mapping is concerned. Existing densities of control in Britain and much of western Europe are already adequate for mapping and cadastral surveys.

Global positioning. The techniques used to establish the positions of reference points within an area to be mapped are similar to those used in navigation. In surveying, however, greater accuracy is required, and this is attainable because the observer and the instrument are stationary on

The imaginary Earth, or geoid

the ground instead of in a ship or aircraft that is not only moving but also subject to accelerations, which make it impossible to use a spirit level for accurate measurements of star elevations.

The technique of locating oneself by observations of celestial objects is rapidly going out of date. In practicing it, the surveyor uses a theodolite with a spirit level to measure accurately the elevations of the Sun at different times of the day or of several known stars in different directions. Each observation defines a line on the Earth's surface on which the observer must be located; several such lines give a fix, the accuracy of which is indicated by how closely these lines meet in a point. For longitude it is necessary also to record the Greenwich Mean Time of each observation. This has been obtained since 1884 by using an accurate chronometer that is checked at least once a day against time signals transmitted telegraphically over land lines and submarine cables or broadcast by radio.

A more recent procedure for global positioning relies on satellites, whose locations at any instant are known precisely because they are being continuously observed from a series of stations in all parts of the world. The coordinates of these stations were established by very large scale triangulation based on a combination of radar observations of distances and measurements of the directions of special balloons or flashing satellites, obtained by photographing them at known instants of time against the background of the fixed stars.

The principal method of using satellites for accurate positioning is based on an application of the Doppler effect. A radio signal is transmitted at a steady frequency by the satellite, but a stationary observer detects a higher frequency as the satellite approaches and a lower one as it recedes. The speed of the frequency drop depends on the distance of the observer from the satellite's track, so a determination of this speed provides a measure of that distance. At the instant of the satellite's closest approach, the observed frequency is the same as that transmitted, so at that time the observer must be located somewhere along the line at right angles to the satellite's track. Since this track over the Earth's surface is accurately known at all times, these data define the observer's position.

Establishing the framework. Most surveying frameworks are erected by measuring the angles and the lengths of the sides of a chain of triangles connecting the points fixed by global positioning. The locations of ground features are then determined in relation to these triangles by less accurate and therefore cheaper methods. Establishing the framework ensures that detail surveys conducted at different times or by different surveyors fit together without overlaps or gaps.

For centuries the corners of these triangles have been located on hilltops, each visible from at least two others, at which the angles between the lines joining them are measured; this process is called triangulation. The lengths of one or two of these lines, called bases, are measured with great care; all the other lengths are derived by trigonometric calculations from them and the angles. Rapid checks on the accuracy are provided by measuring all three angles of each triangle, which must add up to 180 degrees.

In small flat areas, working at large scales, it may be easier to measure the lengths of all the sides, using a tape or a chain, rather than the angles between them; this procedure, called trilateration, was impractical over large or hilly areas until the invention of electromagnetic distance measurement (EDM) in the mid-20th century. This procedure has made it possible to measure distances as accurately and easily as angles, by electronically timing the passage of radiation over the distance to be measured; microwaves, which penetrate atmospheric haze, are used for long distances and light or infrared radiation for short ones. In the devices used for EDM, the radiation is either light (generated by a laser or an electric lamp) or an ultrahigh-frequency radio beam. The light beam requires a clear line of sight; the radio beam can penetrate fog, haze, heavy rain, dust, sandstorms, and some foliage. Both types have a transmitter-receiver at one survey station. At the remote station the light type contains a set of corner mirrors; the high-frequency type incorporates a retrans-

mitter (requiring an operator) identical to the transmitter-receiver at the original station. A corner mirror has the shape of the inside of a corner of a cube; it returns light toward the source from whatever angle it is received, within reasonable limits. A retransmitter must be aimed at the transmitter-receiver.

In both types of instrument, the distance is determined by the length of time it takes the radio or light beam to travel to the target and back. The elapsed time is determined by the shift in phase of a modulating signal superimposed on the carrier beam. Electronic circuitry detects this phase shift and converts it to units of time; the use of more than one modulating frequency eliminates ambiguities that could arise if only a single frequency had been employed.

EDM has greatly simplified an alternative technique, called traversing, for establishing a framework. In traversing, the surveyor measures a succession of distances and the angles between them, usually along a traveled route or a stream. Before EDM was available, traversing was used only in flat or forested areas where triangulation was impossible. Measuring all the distances by tape or chain was tedious and slow, particularly if great accuracy was required, and no check was obtainable until the traverse closed, either on itself or between two points already fixed by triangulation or by astronomical observations.

In both triangulation and traversing, the slope of each measured line must be allowed for so that the map can be reduced to the horizontal and referred to sea level. A measuring tape may be stretched along the ground or suspended between two tripods; in precise work corrections must be applied for the sag, for tension, and for temperature if these differ from the values at which the tape was standardized. In work of the highest order, known as geodetic, the errors must be kept to one millimetre in a kilometre, that is, one part in 1,000,000.

Though for sketch maps the compass or graphic techniques are acceptable for measuring angles, only the theodolite can assure the accuracy required in the framework needed for precise mapping. The theodolite consists of a telescope pivoted around horizontal and vertical axes (Figure 10) so that it can measure both horizontal and vertical angles. These angles are read from circles graduated in degrees and smaller intervals of 10 or 20 minutes. The exact position of the index mark (showing the direction of the line of sight) between two of these graduations is measured on both sides of the circle with the aid of a vernier or a micrometer. The accuracy in modern first-order or geodetic instruments, with five-inch glass circles, is approximately one second of arc, or $1/3,600$ of a degree. With such an instrument a sideways movement of the target of one centimetre can be detected at a distance of two kilometres. By repeating the measurement as many as

Theodolite

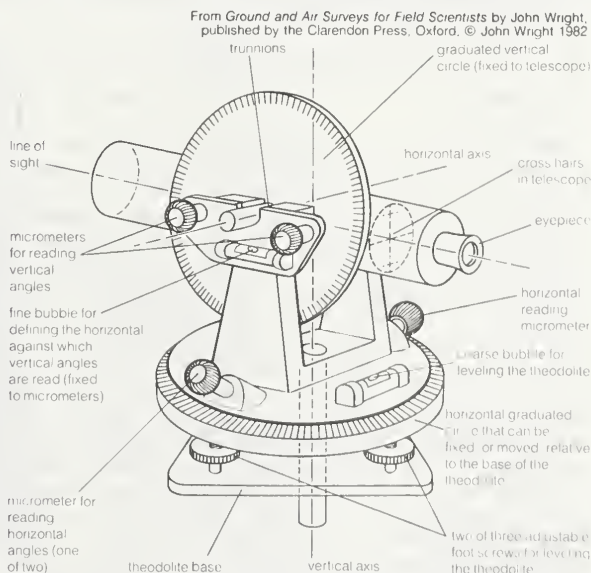


Figure 10: Principal parts of the theodolite. Clamping and slow-motion screws are not shown.

EDM

16 times and averaging the results, horizontal angles can be measured more closely; in geodetic surveying, measurements of all three angles of a triangle are expected to give a sum of 180 degrees within one second of arc.

In the most precise long-distance work, signaling lamps or heliographs reflecting the Sun are used as targets for the theodolite. For less demanding work and work over shorter distances, smaller theodolites with simpler reading systems can be used; targets are commonly striped poles or ranging rods held vertical by an assistant.

An extensive set of these measurements establishes a network of points both on the map, where their positions are plotted by their coordinates, and on the ground, where they are marked by pillars, concrete ground marks, bolts let into the pavement, or wooden pegs of varying degrees of cost and permanence, depending on the importance and accuracy of the framework and the maps to be based on it. Once this framework has been established, the surveyor proceeds to the detail mapping, starting from these ground marks and knowing that their accuracy ensures that the data obtained will fit precisely with similar details obtained elsewhere in the framework.

Detail surveying. The actual depiction of the features to be shown on the map can be performed either on the ground or, since the invention of photography, aviation, and rocketry, by interpretation of aerial photographs and satellite images. On the ground the framework is dissected into even smaller areas as the surveyor moves from one point to another, fixing further points on the features from each position by combinations of angle and distance measurement and finally sketching the features between them freehand. In complicated terrain this operation can be slow and inaccurate, as can be seen by comparing maps made on the ground with those made subsequently from aerial photographs.

Ground survey still has to be used, however, for some purposes; for example, in areas where aerial photographs are hard to get; under the canopy of a forest, where the shape of the ground—not that of the treetops—is required; in very large scale work or close contouring; or if the features to be mapped are not easily identifiable on the aerial photographs, as is the case with property boundaries or zones of transition between different types of soil or vegetation. One of two fundamental differences between ground and air survey is that, as already mentioned, the ground survey interpolates, or sketches, between fixed points, while air survey, using semiautomatic instruments, can trace the features continuously, once the positions of the photographs are known. One effect of this is to show features in uniform detail rather than along short stretches between the points fixed in a ground survey.

The second difference is that in ground survey different techniques and accuracies may be adopted for the horizontal and vertical measurements, the latter usually being more precise. Accurate determinations of heights are required for engineering and planning maps, for example, for railway gradients and particularly for irrigation or drainage networks, since water in open channels does not run uphill.

The methods used for fixing locations within the horizontal detail framework are similar to, but less accurate than, those used for the primary framework. Angles may be measured with a hand-held prismatic compass or graphically with a plane table, or they may be estimated as right angles in the case of points that are offset by short distances from straight lines between points already fixed. Detail points may be located by their distances from two fixed points or by distance and bearing from only one.

The surveyor may record measurements made in the field and plot them there on a sketch board or in the office afterward, but if the country is open and hilly, or even mountainous, the plane table offers the best way of recording the data. A disadvantage of plane-table work is that it cannot be checked in the office, and so it requires greater intelligence and integrity of the surveyor. The plane table reached its most efficient form of use in the Survey of India, begun in 1800, in which large areas were mapped with it by dedicated Indian surveyors. It consists of a flat board that is mounted on a tripod so that it can

be fixed or rotated around a vertical axis. It is set up over a framework point or one end of a measured baseline with its surface (which is covered with paper or other drawing medium) horizontal. It is turned until the line joining its location with another framework point or the other end of the baseline is parallel to the same line as drawn on the paper. This alignment is performed with the aid of an alidade, or sight rule, a straightedge fitted with simple sights. The alidade is then directed toward points on features that are to be fixed, and pencil rays are drawn along the sight rule toward them. The procedure is repeated at the other framework point or the other end of the baseline; the points where the rays intersect on the table will be the map positions of the features (see Figure 11).

From *Ground and Air Surveys for Field Scientists* by John Wright, published by the Clarendon Press, Oxford, © John Wright 1982

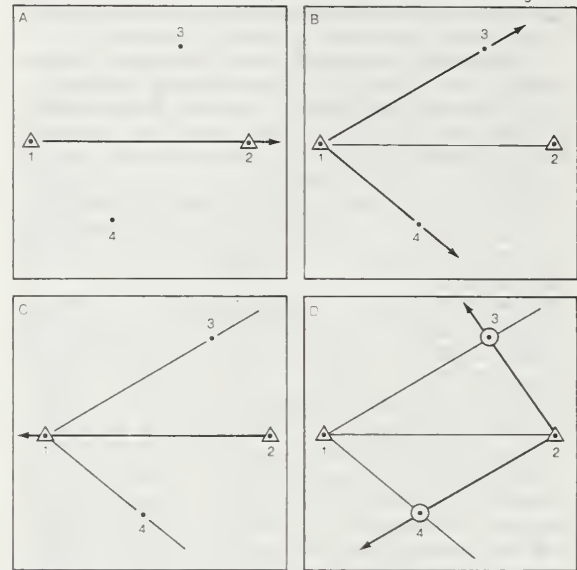


Figure 11: Use of the plane table to determine positions of previously unmapped points (3 and 4) in relation to positions of points of known locations (1 and 2). (A) Table at point 1, rotated to make mapped line 1-2 coincide with actual direction of 2 from 1 and clamped; (B) table at 1, rays drawn in the directions of 3 and 4; (C) table at 2, rotated to make mapped line 2-1 coincide with actual direction of 1 from 2 and clamped; (D) rays drawn in the directions of 3 and 4. The intersections of the rays drawn in steps B and D establish the positions of points 3 and 4.

In surveying for engineering projects, more sophisticated instruments are employed to maximize accuracy. For example, distances may be measured by EDM or by tachymetry, a geometric technique in which the vertical distance on a graduated vertical staff, seen between two stadia hairs in the theodolite eyepiece, is a measure of the horizontal distance between the theodolite and the staff—usually 100 times the difference between the two readings (see Figure 12). This method requires at least one assistant to move the staff from place to place. Modern surveying instruments combine a theodolite, EDM equipment, and a computer that records all the observations and calculates the height differences obtained by measuring vertical angles.

Aerial surveying. Aviation and photography have revolutionized detailed mapping of features visible from the air. An aerial photograph, however, is not a map, as can

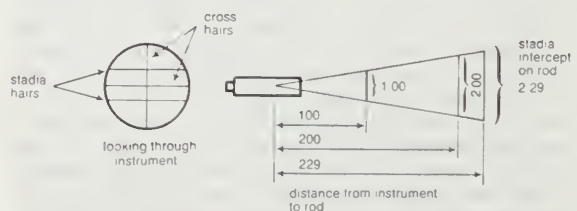


Figure 12: Principles involved in the use of telescope stadia hairs and a staff to determine distance. In the example illustrated, the distance to the rod is equal to 100 times the stadia intercept on the rod (see text).

Use of the plane table

Maps through History



World map with south at the top, by the Arab geographer ash-Sharif al-Idrisi, 1154; from a 1533 copy of al-Idrisi's geography. In the Bodleian Library, University of Oxford, England (MS Pococke 375 folios 3-4).



Hereford map, a T and O map of the world with east at the top and Jerusalem at its centre, by Richard of Haldingham, c. 1280. Actual geography is combined with Christian iconography and mythical elements. Gold and colours on vellum, with oak frame. In the Hereford Cathedral, England.

World map by the Dutch mapmaker Henricus Hondius, 1630. Hondius illustrated the corners of the map with images of his father, Jodocus Hondius (bottom right), Julius Caesar (top left), Claudius Ptolemy (top right), and Gerardus Mercator (bottom left). In the British Library (Maps C.3.d.1).

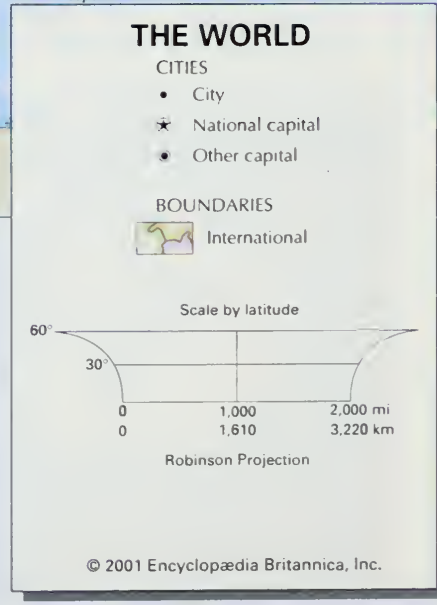


Plate 1: (Top left) Photograph © Woodmansterne; (top right) the Bodleian Library, University of Oxford; (centre and bottom right) by permission of The British Library; (bottom left) courtesy of the John Carter Brown Library at Brown University

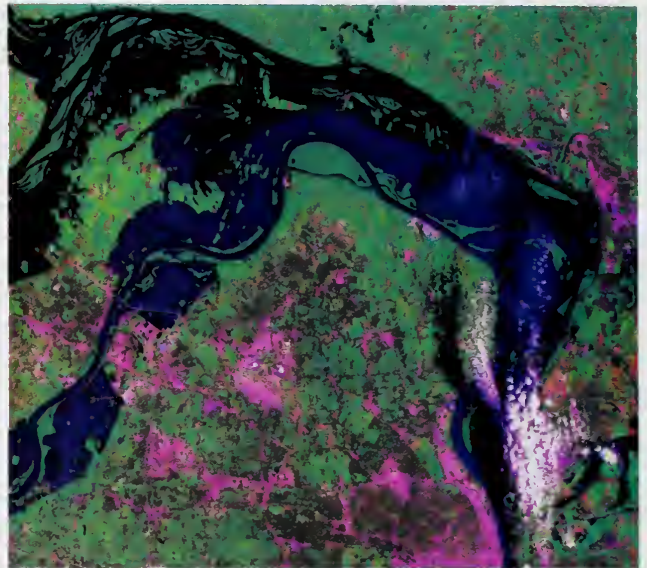


(Left) World map showing the journey of Ferdinand Magellan's fleet across the Pacific, by Giovanni Battista Agnese, 1544. In the John Carter Brown Library, Brown University, Providence, R.I., U.S. (Right) Index chart to the Great Trigonometrical Survey of India, reprinted from *A Memoir on the Indian Surveys*, 2nd edition, by Clements R. Markham (1878). The survey's framework of interlocking triangles took more than 50 years to complete. In the British Library (Maps Ref.K.5).









Landsat 5 Thematic Mapper images of the St. Louis, Mo., U.S., region showing (left) drought conditions on July 4, 1988, and (right) flooding of the Mississippi River (north-south) and the Missouri River (east-west) on July 18, 1993. Vegetation is bright green, bare soils are tan, and urban areas are magenta.



Topographic map of a section of the Wenatchee Mountains, Wash., U.S., with contour lines depicting relief of the landscape.

(Lower left) Axonometric projection map of a section of Manhattan, New York City, produced c. 1997, showing individual buildings in detail, with a three-dimensional effect. Data gathered from site plans, plus aerial and ground photography, were used in creating this map. (Lower right) Aeronautical chart of the southern coast of England. A variety of comprehensive ground information is shown, including emphasis of outstanding landmarks and obstructions, topographic information, principal roads and population centres, visual and radio aids to navigation, controlled airspace, and restricted areas.

Plate 4. Top left and right: Space imaging EOSAT satellite from Washington Atlas & Gazetteer © DeLorme, Vermont, Maine; bottom left © 1997 Ludington Ltd. distributed by Interarts, GeoSystems, Mountain Penn; bottom right: Topographic mapping supplied by Ordnance Survey, the National Mapping Agency of Great Britain © Crown copyright, aeronautical information supplied by the Civil Aviation Authority of Great Britain CAA © copyright





Figure 13: Aerial photograph of (top left) the Houses of Parliament and (right of centre) Westminster Bridge, London. Because they are viewed from a finite altitude, the tops of buildings appear farther from the centre of the picture than the foundations.

By courtesy of J.A. Story & Partners

be seen in Figure 13, a view of the Houses of Parliament and Westminster Bridge, London, from directly overhead. On a map the tops of the towers at the corners of the building should coincide with the corners of the foundations, but in the picture they do not, being displaced radially from the centre. The photograph demonstrates an important property of vertical aerial photographs: angles are correctly represented at their centres, but only there. Similar distortions are present in photographs of hilly ground. This problem may be dealt with in two principal ways, depending on the relative scales of the map and the photographs and on whether contours are required on the map. The older method, adequate for planimetric maps at scales smaller than the photographs, was used extensively during and after World War II to map large areas of desert and thinly populated country; mountainous areas could be sketched in, but the relief was not accurately shown.

As in ground survey, a framework of identified points is necessary before detailed mapping can be carried out from the air. The photographs are ordinarily taken by a vertically aligned camera in a series of strips (see Figure 14) in which each picture overlaps about 60 percent of the preceding one; adjacent strips overlap only slightly. The overlaps make it possible to assemble a low-order framework or control system based on small, recognizable features that appear in more than one photograph. In the

simplest form of this procedure each photograph is replaced by a transparent template on which rays are drawn (or slots are cut) from the centre of the picture to the selected features. The angles between these rays or slots are correct, and slotted templates can be fitted together by inserting studs, which represent the features, into the appropriate slots and sliding the templates so that each

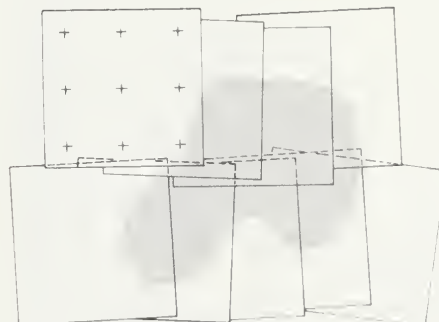


Figure 14: Photogrammetric photographs from two short, overlapping flight strips arranged for supplying mapping details. Photo-control points are shown on only one photograph; shading indicates a typical terrain feature such as a lake (see text).

stud engages the slots in all the pictures showing the corresponding feature. This operation ensures that the centres of the pictures and the selected features are in the correct relationship. The array of overlapping photographs can be expanded or contracted by sliding them about on the work surface as long as the studs remain engaged in the slots, so the assemblage can be positioned, oriented, and scaled by fitting it to at least two—preferably several—ground-control points identified on different photographs.

This technique may be extended by using two additional cameras, one on each side, aimed at right angles to the line of flight and 30 degrees below the horizontal. The photographs taken by the side cameras overlap those taken by the vertical one and also include the horizon; the effect is to widen the strip of ground covered and thus to reduce the amount of flying required. Points in the backgrounds of the oblique photographs can be incorporated in the overlapping array as before to tie the adjacent flight paths together. Photography from high-flying jet aircraft and satellites has rendered this technique obsolete, but before those advances took place it greatly facilitated the mapping of underdeveloped areas.

For the production of maps with accurate contours at scales five or six times that of the photographs, a more sophisticated approach is necessary. The ground-survey effort must be expanded to provide the heights as well as the positions of all the features employed to establish the framework.

From *Ground and Air Surveys for Field Scientists* by John Wright, published by the Clarendon Press, Oxford, © John Wright 1982

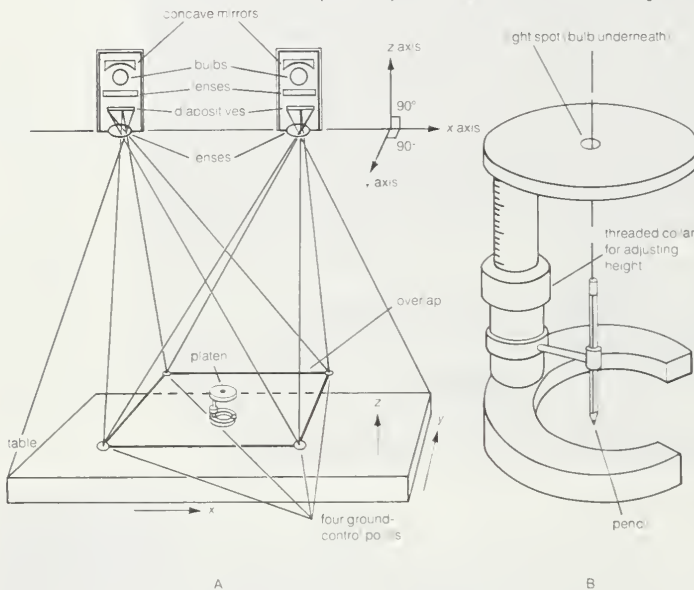


Figure 15: Multiplex projector system of mapping contours and other details from aerial photographs (see text). (A) Arrangement of equipment, showing only two projectors. (B) Details of the platen.

In this technique the details within each segment of the map are based not on individual photographs but on the overlap between two successive ones in the same strip, proceeding from the positions and heights of features in the corners of each area. A three-dimensional model can be created by viewing each pair of consecutive photographs in a stereoscope; by manipulation of a specially designed plotting instrument, the overlapping area can be correctly positioned, scaled, and oriented, and elevations of points within it can be derived from those of the four corner points. These photogrammetric plotting instruments can take several forms. In projection instruments (Figure 15) the photographs are projected onto a table in different colours so that, through spectacles with lenses of complementary colours, each eye sees only one image, and the operator visualizes a three-dimensional model of the ground. A table or platen, with a lighted spot in the middle, can be moved around the model and raised or lowered so that the spot appears to touch the ground while the operator scans any feature, even if it is located on a steep hillside. A pencil directly beneath the spot then plots

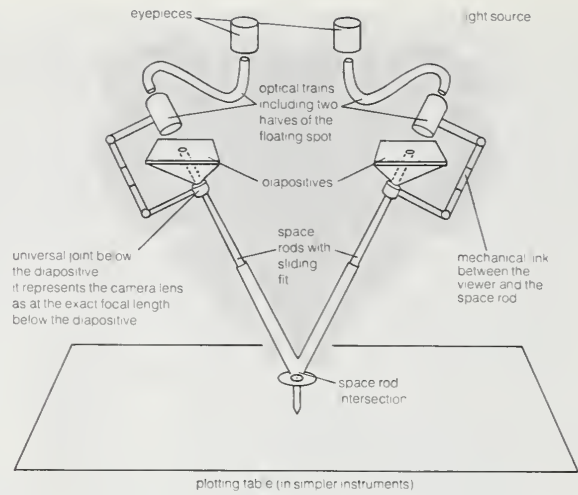


Figure 16: Optical-mechanical plotting instrument (see text). From *Ground and Air Surveys for Field Scientists* by John Wright, published by the Clarendon Press, Oxford, © John Wright 1982

the exact shape and position of the feature on the map. For contouring the platen is fixed at the selected height (at a scale adjusted to that of the model), and the spot is permitted to touch the model surface wherever it will; the pencil then draws the contour.

With more complex mechanical devices, rays of light reaching the aircraft taking the two photographs are represented by rods meeting at a point that represents the position of the feature of the model being viewed (Figure 16). With a complicated system of prisms and lenses the operator, as with projection instruments, sees a spot that can be moved anywhere in the overlap and up or down to touch the model surface, as described above. A mechanical or electronic system moves a pencil into the corresponding position on a plotting table to which the map manuscript is fixed.

With computerized analytic instruments the mechanical operation is limited to measuring coordinates on the two photographs, and the conversion to a three-dimensional model is performed entirely by the computer (Figure 17). It is possible with the most precise plotting instruments of either type to draw a map at four to six times the scale of the photographs and to plot contours accurately at a vertical interval of about one one-thousandth of the height from which the photographs were taken. With such analytic instruments the record can be stored in digital as well as graphic form to be plotted later at any convenient scale.

All these methods produce a line or drawn map; some of them also create a data file on disk or tape, containing the coordinates of all the lines and other features on the map. On the other hand, aerial photographs can be combined and printed directly to form a photomap. For flat areas this operation requires simply cutting and pasting the photographs together into a mosaic. For greater accuracy the centres of the photographs may be aligned by the use of slotted templates as described above to produce a photomap called a controlled mosaic.

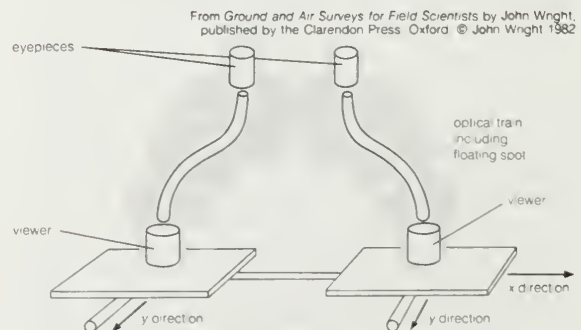


Figure 17: Geometry of analytical and approximate photogrammetric instruments. Both diapositives (or paper prints) are in the same plane and can move only in the x and y directions together and relative to each other.

A much more precise technique is based on the use of an orthophotoscope. With this device, overlapping photographs are employed just as in the stereoscopic plotter already described, but the instrument, rather than the manual tracing of the features and contours, scans the overlap and produces an orthophotograph by dividing the area into small sections, each of which is correctly scaled. This procedure is best applied to areas of low relief without tall buildings; the resulting maps can then be substituted for line maps in rural areas where they are particularly useful in planning resettlement in agricultural projects. Because no fair drawing is required, the final printed map can be produced much more quickly and cheaply than would otherwise be possible.

Hydrography. Surveying of underwater features, or hydrographic surveying, formerly required techniques very different from ground surveying, for two reasons: the surveyor ordinarily was moving instead of stationary, and the surface being mapped could not be seen. The first problem, making it difficult to establish a framework except near land or in shoal areas, was dealt with by dead reckoning between points established by astronomical fixes. In effect a traverse would be run with the ship's bearing measured by compass and distances obtained either by measuring speed and time or by a modern log that directly records distances. These have to be checked frequently, because however accurate the log or airspeed indicator and compass, the track of a ship or aircraft is not the same as its course. Crosscurrents or winds continually drive the craft off course, and those along the course affect the speed and the distance run over the ground beneath.

The only way a hydrographer could chart the seabed before underwater echo sounding and television became available was to cast overboard at intervals a sounding line with a lead weight at the end and measure the length of the line paid out when the weight hit the bottom. The line was marked in fathoms, that is, units of one one-thousandth of a nautical mile, or approximately six feet (1.8 metres).

Sounding by lead line is obviously very slow, especially in deep waters, and the introduction of echo sounding in the early 20th century marked a great improvement. It was made possible by the invention of electronic devices for the measurement of short intervals of time. Echo sounding depends on timing the lapse between the transmission of a short loud noise or pulse and its return from the target—in this case the bottom of the sea or lake. Sound travels about 5,000 feet (1,500 metres) per second in water, so that an accuracy of a few milliseconds in measurements of the time intervals gives depths within a few feet.

The temperature and density of water affect the speed at which sound waves travel through it, and allowances have to be made for variations in these properties. The reflected signals are recorded several times a second on a moving strip of paper, showing to scale the depth beneath the ship's track. The echoes may also show other objects, such as schools of fish, or they may reveal the dual nature of the bottom, where a layer of soft mud may overlie rock. Originally only the depth that was directly beneath the ship was measured, leaving gaps between the ship's tracks. Later inventions, which include sideways-directed sonar and television cameras, have made it possible to fill these gaps. While measurements of depths away from the ship's track are not so accurate, the pictures reveal any dangerous objects such as rock pinnacles or wrecks, and the survey vessel can then be diverted to survey them in detail.

Modern position-fixing techniques using radar have made the whole process much simpler, for the ship's location is now known continuously with reference to fixed stations on shore or to satellite tracks. Another modern technique is the use of pictures taken from aircraft or satellites to indicate the presence and shape of shoal areas and to aid the planning of their detailed survey.

An alternative to the use of radar or satellite signals for continuous and automatic recording of a ship's position is the employment of inertial guidance systems. These devices, developed to satisfy military requirements, detect every acceleration involved in the motion of a craft from its known starting point and convert them and the elapsed

time into a continuous record of the distance and direction traveled.

For studying the seabed in detail, the bottom of the sounding lead was hollowed to hold a charge of grease to pick up a sample from the sea floor. Today television cameras can be lowered to transmit pictures back to the survey ship, though their range is limited by the extent to which light can penetrate the water, which often is murky. Ordinary cameras also are used in pairs for making stereoscopic pictures of underwater structures such as drilling rigs or the wreckage of ancient ships.

Height determination. Heights of surface features above sea level are determined in four main ways: by spirit leveling, by measuring vertical angles and distances, by measuring differences in atmospheric pressure, and, since the late 20th century, by using three-dimensional satellite or inertial systems. Of these the first is the most accurate; the second is next in accuracy but faster; the third is least accurate but can be fastest if heights are to be measured at well-separated points. The last two techniques require sophisticated equipment that is still very expensive.

In spirit leveling the surveyor has for centuries used a surveying level, which consists of a horizontal telescope fitted with cross hairs, rotating around a vertical axis on a tripod, with a very sensitive spirit level fixed to it; the instrument is adjusted until the bubble is exactly centred. The reading on a graduated vertical staff is observed through the telescope. If such staffs are placed on successive ground points, and the telescope is truly level, the difference between the readings at the cross hairs will equal that between the heights of the points. By moving the level and the staffs alternately along a path or road and repeating this procedure, differences in height can be accurately measured over long horizontal distances.

In the most precise work, over a distance of 100 kilometres the error may be kept to less than a centimetre. To achieve this accuracy great care has to be taken. The instrument must have a high-magnification telescope and a very sensitive bubble, and the graduated scale on the staff must be made of a strip of Invar (an alloy with a very small coefficient of thermal expansion). Moreover, the staffs must be placed on pegs or special heavy steel plates, and the distance between them and the level must always be the same to cancel the effects of aerial refraction of the light.

In less precise work a single wooden staff can be used; for detailed leveling of a small area, the staff is moved from one point to another without moving the level so that heights can be measured within a radius of about 100 metres (Figure 18). The distances of these points from the instrument can be measured by tape or, more commonly, by recording not only the reading at the central cross hair in the field of view of the telescope but also those at the stadia hairs, that is by tachymetry, as described above. The bearing of each point is observed by compass or on the horizontal circle of the level so that it can be plotted or drawn on the map.

From *Ground and Air Surveys for Field Scientists* by John Wright, published by the Clarendon Press, Oxford, © John Wright 1962

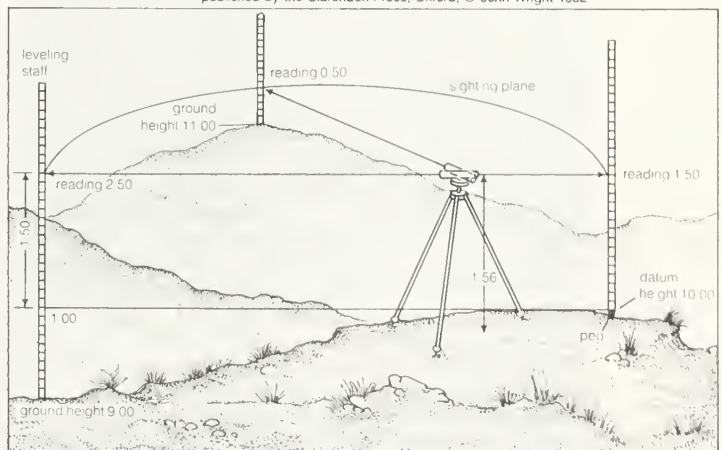


Figure 18: Spirit leveling to measure heights of ground points relative to that of a datum point (see text).

Spirit leveling

Echo sounding

Since the 1950s levels have been introduced in which the line of sight is automatically leveled by passage through a system of prisms in a pendulum, thus removing the need to check the bubble. The disadvantage of spirit leveling is the large number of times the instrument has to be moved and realigned, particularly on steep hills; it is used primarily along practically flat stretches of ground.

For faster work in hilly areas, where lower accuracies usually are acceptable, trigonometric height determination is employed using a theodolite to measure vertical angles and measuring or calculating the distances by triangulation. This procedure is particularly useful in obtaining heights throughout a major framework of triangulation or traverse where most of the points are on hilltops. To increase precision, the observations are made simultaneously in both directions so that aerial refraction is eliminated; this is done preferably around noon, when the air is well mixed.

The third method of height determination depends on measurements of atmospheric pressure differences with a sensitive aneroid barometer, which can respond to pressure differences small enough to correspond to a foot or two (0.3 to 0.6 metre) in height. The air pressure changes constantly, however, and to obtain reliable results it is necessary to use at least two barometers; one at a reference point of known height is read at regular intervals while the surveyor proceeds throughout the area, recording locations, times, and barometer readings. Comparison of readings made at the same time then gives the height differences.

An alternative to the barometer for pressure measurement is an apparatus for measuring the boiling point of a liquid, because this temperature depends on the atmospheric pressure. Early explorers determined heights in this way, but the results were very rough; this technique was not accurate enough for surveyors until sensitive methods for temperature measurement were developed. The airborne profile recorder is a combination of this refined apparatus with a radar altimeter to measure the distance to the ground below an aircraft.

Analysis of the signals received simultaneously from several satellites gives heights as accurately as positions. Heights determined in this way are useful in previously unmapped areas as a check on results obtained by faster relative methods, but they are not accurate enough for mapping developed areas or for engineering projects. All-terrain vehicles or helicopters can carry inertial systems accurate enough to provide approximate heights suitable for aerial surveys of large areas within a framework of points established more accurately by spirit leveling.

(J.W.Wr./Ed.)

BIBLIOGRAPHY

History of mapping and surveying: A list of sources on the topic is presented in WALTER W. RISTOW, *Guide to the History of Cartography: An Annotated List of References on the History of Maps and Mapmaking* (1973). ROBERT C. DURU, *Maps and Map Making* (1977); and G.R. CRONE, *Maps and Their Makers: An Introduction to the History of Cartography*, 5th ed. (1978), are brief overviews. Essential aspects are explored in LEO BAGROW, *History of Cartography*, 2nd ed., rev. by R.A. SKELTON (1985, originally published in German, 1951); CHARLES BRICKER, *Landmarks of Mapmaking: An Illustrated Survey of Maps and Mapmakers* (1968; also published as *A History of Cartography: 2500 Years of Maps and Mapmakers*, 1969; reissued 1977); LLOYD A. BROWN, *The Story of Maps* (1949, reprinted 1979), and *Map Making: The Art That Became a Science* (1960); EDWARD LYNAM, *The Mapmaker's Art: Essays on the History of Maps* (1953); and JOHN NOBLE WILFORD, *The Mapmakers* (1981). Maps of specific periods are studied in CHARLES H. HAPGOOD, *Maps of the Ancient Sea Kings: Evidence of Advanced Civilization in the Ice Age*, rev. ed. (1979); R.A. SKELTON, *Decorative Printed Maps of the 15th to 18th Centuries* (1952, reprinted 1966); HUGH CORTAZZI, *Isles of Gold: Antique*

Maps of Japan (1983); and RAYMOND LISTER, *How to Identify Old Maps and Globes: With a List of Cartographers, Engravers, Publishers, and Printers Concerned with Printed Maps and Globes from c. 1500 to c. 1850* (1965). The development of cartographic institutions can be traced in MARY BLEWITT, *Surveys of the Seas: A Brief History of British Hydrography* (1957); SIR ARCHIBALD DAY, *The Admiralty Hydrographic Service, 1795-1919* (1967); ADRIAN H.W. ROBINSON, *Marine Cartography in Britain: A History of the Sea Chart to 1855* (1962); and G.S. RITCHIE, *The Admiralty Chart: British Naval Hydrography in the Nineteenth Century* (1967).

Mapmaking: The procedures involved in creating a map are described in many books, including WELLMAN CHAMBERLIN, *The Round Earth on Flat Paper: Map Projections Used by Cartographers* (1947); ARTHUR D. MERRIMAN, *An Introduction to Map Projections* (1947); ERWIN J. RAISZ, *General Cartography*, 2nd ed. (1948), and *Principles of Cartography* (1962); EDUARD IMHOF, *Cartographic Relief Presentation* (1982; originally published in German, 1965); J.S. KEATES, *Cartographic Design and Production* (1973); ARTHUR H. ROBINSON et al., *Elements of Cartography*, 5th ed. (1984); JOHN LOXTON, *Practical Map Production* (1980); and MORRIS M. THOMPSON, *Maps for America: Cartographic Products of the U.S. Geological Survey and Others*, 2nd ed. (1981). A range of special projects, of which cartography is an essential part, are described in A.H.A. HOGG, *Surveying for Archaeologists and Other Fieldworkers* (1980); TEODOR J. BLACHUT, ADAM CHRZANOWSKI, and JOUKO H. SAASTAMOINEN, *Urban Surveying and Mapping* (1979); P.F. DALE, *Cadastral Surveys Within the Commonwealth* (1976); S. ROWTON SIMPSON, *Land Law and Registration* (1976); R.A. SKELTON, *The Legal Elements of Boundaries and Adjacent Properties* (1930); CURTIS M. BROWN, WALTER G. ROBILLARD, and DONALD A. WILSON, *Evidence and Procedures for Boundary Location*, 2nd ed. (1981); T.W. BIRCH, *Maps: Topographical and Statistical*, 2nd ed. (1964, reprinted with corrections, 1976); B.W. LUCKE, *A Course on the Chart* (1966); M. CHRIS and G.R. HAYES, *An Introduction to Charts and Their Use*, 4th ed. (1977); and D.A. MOORE, *Marine Chartwork*, 2nd ed. (1981).

Technical developments: The many special studies in the field include CHESTER C. SLAMA (ed.), *Manual of Photogrammetry*, 4th ed. (1980); JOHN WRIGHT, *Ground and Air Survey for Field Scientists* (1982); C.D. BURNSIDE, *Electromagnetic Distance Measurement*, 2nd ed. (1982); and WILLIAM RITCHIE et al., *Mapping for Field Scientists: A Problem-Solving Approach* (1977). Current developments are covered in special journals: *Surveying and Mapping* (quarterly); *Photogrammetria* (bimonthly); *Military Engineer* (bimonthly); *The American Cartographer* (semiannual); *Cartography* (semiannual); and *Cartographica* (quarterly).

Modern surveying: UNITED STATES GOVERNMENT PRINTING OFFICE, *Surveying and Mapping* (1982), is a subject bibliography listing current sources. Informative introductory texts include WILLIAM C. WATTLES, *Land Survey Descriptions*, rev. ed., edited by GURDON H. WATTLES (1974); CHARLES A. HERUBIN, *Principles of Surveying*, 3rd ed. (1982); J.G. OLLIVER and J. CLENDINNING, *Principles of Surveying: Plane Surveying*, 4th ed. (1978); JERRY A. NATHANSON and PHILIP KISSAM, *Surveying Practice*, 4th ed. (1988); R.H. DUGDALE, *Surveying*, 3rd ed. (1980); and JACK C. MCCORMAC, *Surveying Fundamentals* (1983). Comprehensive treatments include CHARLES B. BREED and GEORGE L. HOSMER, *The Principles and Practices of Surveying*, 11th ed., rev. by W. FAIG, 2 vol. (1977); RAYMOND E. DAVIS et al., *Surveying, Theory and Practice*, 6th ed. (1981); RUSSELL C. BRINKER and PAUL R. WOLF, *Elementary Surveying*, 7th ed. (1984); FRANCIS H. MOFFITT and HARRY BOUCHARD, *Surveying*, 7th ed. (1982); and G. BOMFORD, *Geodesy*, 4th ed. (1980). Instruments used in the practice are discussed in J. CLENDINNING and J.G. OLLIVER, *Principles and Use of Surveying Instruments*, 3rd ed. (1969, reissued 1972); and M.A.R. COOPER, *Modern Theodolites and Levels*, 2nd ed. (1982). Hydrography is the subject of A.E. INGHAM (ed.), *Sea Surveying*, 2 vol. (1975); and MELVIN J. UMBACH, *Hydrographic Manual*, 4th ed. (1976). INTERNATIONAL HYDROGRAPHIC BUREAU, *Hydrographic Dictionary*, 3rd ed. (1970), contains useful information; and *International Hydrographic Bulletin* (monthly) and *International Hydrographic Review* (semiannual), published by the same organization, provide data on current research and developments.

(C.F.F./J.W.Wr.)

Marketing and Merchandising

Marketing is a process whose principal function is to promote and facilitate exchange. Through marketing, individuals and groups obtain what they need and want by exchanging products and services with other parties. Such a process can occur only when there are at least two parties, each of whom has something to offer. In addition, exchange cannot occur unless the parties are able to communicate about and to deliver what they offer. Marketing is not a coercive process: all parties must be free to accept or reject what others are offering. So defined, marketing is distinguished from other modes of obtaining desired goods, such as through self-production, begging, theft, or force.

Marketing is not confined to any particular type of economy, because goods must be exchanged and therefore marketed in all economies and societies except perhaps in

the most primitive. Furthermore, marketing is not a function that is limited to profit-oriented business; even such institutions as hospitals, schools, and museums engage in some forms of marketing. Within the broad scope of marketing, merchandising is concerned more specifically with promoting the sale of goods and services to consumers (*i.e.*, retailing) and hence is more characteristic of free-market economies.

Based on these criteria, marketing can take a variety of forms: it can be a set of functions, a department within an organization, a managerial process, a managerial philosophy, and a social process.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, sections 532 and 533, and the *Index*.

This article is divided into the following sections:

-
- The evolving discipline of marketing 495
 - Roles of marketing 496
 - The marketing process 496
 - Strategic marketing analysis 496
 - Marketing-mix planning 496
 - Product
 - Price
 - Place
 - Promotion
 - Marketing implementation 498
 - Marketing evaluation and control 498
 - Customers 498
 - Consumer customers 499
 - Factors influencing consumers
 - Consumer buying tasks
 - The consumer buying process
 - Business customers 500
 - Factors influencing business customers
 - The business buying process
 - Marketing intermediaries 500
 - Channel functions and flows
 - Management of channel systems
 - Merchandise distributors 501
 - Wholesalers 501
 - Merchant wholesalers
 - Brokers and agents
 - Manufacturers' and retailers' branches and offices
 - Retailers 502
 - The history of retailing
 - Store retailers
 - Specialty stores
 - Department stores
 - Supermarkets
 - Convenience stores
 - Superstores
 - Discount stores
 - Off-price retailers
 - Nonstore retailers
 - Direct selling
 - Direct marketing
 - Automatic vending
 - Retail organizations
 - Corporate chains
 - Voluntary chains and retailer cooperatives
 - Consumer cooperatives
 - Franchise organizations
 - Merchandising conglomerates
 - Marketing facilitators 504
 - Advertising agencies 504
 - Market research firms 504
 - Transportation firms 505
 - Warehousing firms 505
 - Marketing in different sectors 505
 - The government market
 - Consumer-goods marketing
 - Services marketing
 - Business marketing
 - Nonprofit marketing
 - Social marketing
 - Place marketing
 - Economic and social aspects of marketing 507
 - Marketing and individual welfare
 - Marketing and societal welfare
 - Marketing's contributions to individuals and society
 - Bibliography 507
-

The evolving discipline of marketing

The marketing discipline had its origins in the early 20th century as an offspring of economics. Economic science had neglected the role of middlemen and the role of functions other than price in the determination of demand levels and characteristics. Early marketing economists examined agricultural and industrial markets and described them in greater detail than the classical economists. This examination resulted in the development of three approaches to the analysis of marketing activity: the commodity, the institution, and the function.

Commodity analysis studies the ways in which a product or product group is brought to market. A commodity analysis of milk, for example, traces the ways in which milk is collected at individual dairy farms, transported to and processed at local dairy cooperatives, and shipped to grocers and supermarkets for consumer purchase. Institutional analysis describes the types of businesses that play a prevalent role in marketing, such as wholesale or retail institutions. For instance, an institutional analysis

of clothing wholesalers examines the ongoing concerns that wholesalers face in order to ensure both the correct supply for their customers and the appropriate inventory and shipping capabilities. Finally, a functional analysis examines the general tasks that marketing performs. For example, any marketing effort must ensure that the product is transported from the supplier to the customer. In some industries, this transportation function may be handled by a truck, while in others it may be done by mail, facsimile, television signal, or airline. All these institutions perform the same function.

As the study of marketing became more prevalent throughout the 20th century, large companies—particularly mass consumer manufacturers—began to recognize the importance of market research, better product design, effective distribution, and sustained communication with consumers in the success of their brands. Marketing concepts and techniques later moved into the industrial-goods sector and subsequently into the services sector. It soon became apparent that organizations and individuals market not only goods and services but also ideas (so-

Rise of
marketing

cial marketing), places (location marketing), personalities (celebrity marketing), events (event marketing), and even the organizations themselves (public relations).

Roles of marketing

As marketing developed, it took a variety of forms. It was noted above that marketing can be viewed as a set of functions in the sense that certain activities are traditionally associated with the exchange process. A common but incorrect view is that selling and advertising are the only marketing activities. Yet, in addition to promotion, marketing includes a much broader set of functions, including product development, packaging, pricing, distribution, and customer service.

Many organizations and businesses assign responsibility for these marketing functions to a specific group of individuals within the organization. In this respect, marketing is a unique and separate entity. Those who make up the marketing department may include brand and product managers, marketing researchers, sales representatives, advertising and promotion managers, pricing specialists, and customer service personnel.

As a managerial process, marketing is the way in which an organization determines its best opportunities in the marketplace, given its objectives and resources. The marketing process is divided into a strategic and a tactical phase. The strategic phase has three components—segmentation, targeting, and positioning (STP). The organization must distinguish among different groups of customers in the market (segmentation), choose which group(s) it can serve effectively (targeting), and communicate the central benefit it offers to that group (positioning). The marketing process includes designing and implementing various tactics, commonly referred to as the “marketing mix,” or the “4 Ps”: product, price, place (or distribution), and promotion. The marketing mix is followed by evaluating, controlling, and revising the marketing process to achieve the organization’s objectives (see below the section *Marketing-mix planning*).

The managerial philosophy of marketing puts central emphasis on customer satisfaction as the means for gaining and keeping loyal customers. Marketers urge their organizations to carefully and continually gauge target customers’ expectations and to consistently meet or exceed these expectations. In order to accomplish this, everyone in all areas of the organization must focus on understanding and serving customers; it will not succeed if all marketing occurs only in the marketing department. Marketing, consequently, is far too important to be done solely by the marketing department. Marketers also want their organizations to move from practicing transaction-oriented marketing, which focuses on individual exchanges, to relationship-driven marketing, which emphasizes serving the customer over the long term. Simply getting new customers and losing old ones will not help the organization achieve its objectives.

Finally, marketing is a social process that occurs in all economies, regardless of their political structure and orientation. It is the process by which a society organizes and distributes its resources to meet the material needs of its citizens. However, marketing activity is more pronounced under conditions of goods surpluses than goods shortages. When goods are in short supply, consumers are usually so desirous of goods that the exchange process does not require significant promotion or facilitation. In contrast, when there are more goods and services than consumers need or want, companies must work harder to convince customers to exchange with them.

The marketing process

The marketing process consists of four elements: strategic marketing analysis, marketing-mix planning, marketing implementation, and marketing control.

STRATEGIC MARKETING ANALYSIS

The aim of marketing in profit-oriented organizations is to meet needs profitably. Companies must therefore first

define which needs—and whose needs—they can satisfy. For example, the personal transportation market consists of people who put different values on an automobile’s cost, speed, safety, status, and styling. No single automobile can satisfy all these needs in a superior fashion; compromises have to be made. Furthermore, some individuals may wish to meet their personal transportation needs with something other than an automobile, such as a motorcycle, a bicycle, or a bus or other form of public transportation. Because of such variables, an automobile company must identify the different preference groups, or segments, of customers and decide which group(s) they can target profitably.

Segments can be divided into even smaller groups, called subsegments or niches. A niche is defined as a small target group that has special requirements. For example, a bank may specialize in serving the investment needs of not only senior citizens but also senior citizens with high incomes and perhaps even those with particular investment preferences. It is more likely that larger organizations will serve the larger market segments (mass marketing) and ignore niches. As a result, smaller companies typically emerge that are intimately familiar with a particular niche and specialize in serving its needs.

A growing number of companies are now trying to serve “segments of one.” They attempt to adapt their offer and communication to each individual customer. This is understandable, for instance, with large industrial companies that have only a few major customers. For example, The Boeing Company (United States) designs its 747 planes differently for each major customer, such as United Airlines, Inc., or American Airlines, Inc. Serving individual customers is increasingly possible with the advent of database marketing, through which individual customer characteristics and purchase histories are retained in company information systems. Even mass-marketing companies, particularly large retailers and catalog houses, compile comprehensive data on individual customers and are able to customize their offerings and communications.

A key step in marketing strategy, known as positioning, involves creating and communicating a message that clearly establishes the company or brand in relation to competitors. Thus, Volvo Aktiebolaget (Sweden) has positioned its automobile as the “safest,” and Daimler-Benz AG (Germany), manufacturer of Mercedes-Benz vehicles, has positioned its car as the best “engineered.” Some products may be positioned as “outstanding” in two or more ways. However, claiming superiority along several dimensions may hurt a company’s credibility because consumers will not believe that any one offering can excel in all dimensions. Furthermore, although the company may communicate a particular position, customers may perceive a different image of the company as a result of their actual experiences with the company’s product or through word of mouth.

MARKETING-MIX PLANNING

Having developed a strategy, a company must then decide which tactics will be most effective in achieving strategy goals. Tactical marketing involves creating a marketing mix of four components—product, price, place, promotion—that fulfills the strategy for the targeted set of customer needs.

Product. The first marketing-mix element is the product, which refers to the offering or group of offerings that will be made available to customers. In the case of a physical product, such as a car, a company will gather information about the features and benefits desired by a target market. Before assembling a product, the marketer’s role is to communicate customer desires to the engineers who design the product or service. This is in contrast to past practice, when engineers designed a product based on their own preferences, interests, or expertise and then expected marketers to find as many customers as possible to buy this product. Contemporary thinking calls for products to be designed based on customer input and not solely on engineers’ ideas.

In traditional economies, the goods produced and consumed often remain the same from one generation

Relation to supply and demand

Consumer input in product development

to the next—including food, clothing, and housing. As economies develop, the range of products available tends to expand, and the products themselves change. In contemporary industrialized societies, products, like people, go through life cycles: birth, growth, maturity, and decline. This constant replacement of existing products with new or altered products has significant consequences for professional marketers. The development of new products involves all aspects of a business—production, finance, research and development, and even personnel administration and public relations.

Packaging and branding are also substantial components in the marketing of a product. Packaging in some instances may be as simple as customers in France carrying long loaves of unwrapped bread or small produce dealers in Italy wrapping vegetables in newspapers or placing them in customers' string bags. In most industrialized countries, however, the packaging of merchandise has become a major part of the selling effort, as marketers now specify exactly the types of packaging that will be most appealing to prospective customers. The importance of packaging in the distribution of the product has increased with the spread of self-service purchases—in wholesaling as well as in retailing. Packaging is sometimes designed to facilitate the use of the product, as with aerosol containers for room deodorants. In Europe such condiments as mustard, mayonnaise, and ketchup are often packaged in tubes. Some packages are reusable, making them attractive to customers in poorer countries where metal containers, for instance, are often highly prized.

The same general marketing approach about the product applies to the development of service offerings as well. For example, a health maintenance organization (HMO) must design a contract for its members that describes which medical procedures will be covered, how much physician choice will be available, how out-of-town medical costs will be handled, and so forth. In creating a successful service mix, the HMO must choose features that are preferred and expected by target customers, or the service will not be valued in the marketplace.

Price. The second marketing-mix element is price. Ordinarily companies determine a price by gauging the quality or performance level of the offer and then selecting a price that reflects how the market values its level of quality. However, marketers also are aware that price can send a message to a customer about the product's presumed quality level. A Mercedes-Benz vehicle is generally considered to be a high-quality automobile, and it therefore can command a high price in the marketplace. But, even if the manufacturer could price its cars competitively with economy cars, it might not do so, knowing that the lower price might communicate lower quality. On the other hand, in order to gain market share, some companies have moved to "more for the same" or "the same for less" pricing, which means offering prices that are consistently lower than those of their competitors. This kind of discount pricing has caused firms in such industries as airlines and pharmaceuticals (which used to charge a price premium based on their past brand strength and reputation) to significantly reevaluate their marketing strategies.

Place. Place, or where the product is made available, is the third element of the marketing mix and is most commonly referred to as distribution. When a product moves along its path from producer to consumer, it is said to be following a channel of distribution. For example, the channel of distribution for many food products includes food-processing plants, warehouses, wholesalers, and supermarkets. By using this channel, a food manufacturer makes its products easily accessible by ensuring that they are in stores that are frequented by those in the target market. In another example, a mutual funds organization makes its investment products available by enlisting the assistance of brokerage houses and banks, which in turn establish relationships with particular customers. However, each channel participant can handle only a certain number of products: space at supermarkets is limited, and investment brokers can keep abreast of only a limited number of mutual funds. Because of this, some marketers may decide to skip steps in the channel and instead mar-

ket directly to buyers through direct mail, telemarketing, door-to-door selling, shopping via television (a growing trend in the late 20th century), or factory outlets.

Promotion. Promotion, the fourth marketing-mix element, consists of several methods of communicating with and influencing customers. The major tools are sales force, advertising, sales promotion, and public relations.

Sales representatives are the most expensive means of promotion, because they require income, expenses, and supplementary benefits. Their ability to personalize the promotion process makes salespeople most effective at selling complex goods, big-ticket items, and highly personal goods—for example, those related to religion or insurance. Salespeople are trained to make presentations, answer objections, gain commitments to purchase, and manage account growth. Some companies have successfully reduced their sales-force costs by replacing certain functions (for example, finding new customers) with less expensive methods (such as direct mail and telemarketing).

Advertising includes all forms of paid, nonpersonal communication and promotion of products, services, or ideas by a specified sponsor. Advertising appears in such media as print (newspapers, magazines, billboards, flyers) or broadcast (radio, television). Print advertisements typically consist of a picture, a headline, information about the product, and occasionally a response coupon. Broadcast advertisements consist of an audio or video narrative that can range from short 15-second spots to longer segments known as infomercials, which generally last 30 or 60 minutes.

While advertising presents a reason to buy a product, sales promotion offers a short-term incentive to purchase. Sales promotions often attract brand switchers (those who are not loyal to a specific brand) who are looking primarily for low price and good value. Thus, especially in markets where brands are highly similar, sales promotions can cause a short-term increase in sales but little permanent gain in market share. Alternatively, in markets where brands are quite dissimilar, sales promotions can alter market shares more permanently. The use of promotions has risen considerably during the late 20th century. This is due to a number of factors within companies, including an increased sophistication in sales promotion techniques and greater pressure to increase sales. Several market factors also have fostered this increase, including a rise in the number of brands (especially similar ones) and a decrease in the efficiency of traditional advertising due to increasingly fractionated consumer markets.

Public relations, in contrast to advertising and sales promotion, generally involves less commercialized modes of communication. Its primary purpose is to disseminate information and opinion to groups and individuals who have an actual or potential impact on a company's ability to achieve its objectives. In addition, public relations specialists are responsible for monitoring these individuals and groups and for maintaining good relationships with them. One of their key activities is to work with news and information media to ensure appropriate coverage of the company's activities and products. Public relations specialists create publicity by arranging press conferences, contests, meetings, and other events that will draw attention to a company's products or services. Another public relations responsibility is crisis management—that is, handling situations in which public awareness of a particular issue may dramatically and negatively impact the company's ability to achieve its goals. For example, when it was discovered that some bottles of Perrier sparkling water might have been tainted by a harmful chemical, Source Perrier, SA's public relations team had to ensure that the general consuming public did not thereafter automatically associate Perrier with tainted water. Other public relations activities include lobbying, advising management about public issues, and planning community events.

Because public relations does not always seek to impact sales or profitability directly, it is sometimes seen as serving a function that is separate from marketing. However, some companies recognize that public relations can work in conjunction with other marketing activities to facilitate the exchange process directly and indirectly. These

Types of
advertising

organizations have established marketing public relations departments to directly support corporate and product promotion and image management.

(P.Ko./K.A.G./Jo.D.H.)

MARKETING IMPLEMENTATION

Companies have typically hired different agencies to help in the development of advertising, sales promotion, and publicity ideas. However, this often results in a lack of coordination between elements of the promotion mix. When components of the mix are not all in harmony, a confusing message may be sent to consumers. For example, a print advertisement for an automobile may emphasize the car's exclusivity and luxury, while a television advertisement may stress rebates and sales, clashing with this image of exclusivity. Alternatively, by integrating the marketing elements, a company can more efficiently utilize its resources. Instead of individually managing four or five different promotion processes, the company manages only one. In addition, promotion expenditures are likely to be better allocated, because differences among promotion tools become more explicit. This reasoning has led to integrated marketing communications, in which all promotional tools are considered to be part of the same effort, and each tool receives full consideration in terms of its cost and effectiveness.

MARKETING EVALUATION AND CONTROL

No marketing process, even the most carefully developed, is guaranteed to result in maximum benefit for a company. In addition, because every market is changing constantly, a strategy that is effective today may not be effective in the future. It is important to evaluate a marketing program periodically to be sure that it is achieving its objectives. There are four types of marketing control, each of which has a different purpose: annual-plan control, profitability control, efficiency control, and strategic control.

The basis of annual-plan control is managerial objectives—that is to say, specific goals, such as sales and profitability, that are established on a monthly or quarterly basis. Organizations use five tools to monitor plan performance. The first is sales analysis, in which sales goals are compared with actual sales and discrepancies are explained or accounted for. A second tool is market-share analysis, which compares a company's sales with those of its competitors. Companies can express their market share in a number of ways, by comparing their own sales to total market sales, sales within the market segment, or sales of the segment's top competitors. Third, marketing expense-to-sales analysis gauges how much a company spends to achieve its sales goals. The ratio of marketing expenses to sales is expected to fluctuate, and companies usually establish an acceptable range for this ratio. In contrast, financial analysis estimates such expenses (along with others) from a corporate perspective. This includes a comparison of profits to sales (profit margin), sales to assets (asset turnover), profits to assets (return on assets), assets to worth (financial leverage), and, finally, profits to worth (return on net worth). Finally, companies measure customer satisfaction as a means of tracking goal achievement. Analyses of this kind are generally less quantitative than those described above and may include complaint and suggestion systems, customer satisfaction surveys, and careful analysis of reasons why customers switch to a competitor's product.

Profitability control and efficiency control allow a company to closely monitor its sales, profits, and expenditures. Profitability control demonstrates the relative profit-earning capacity of a company's different products and consumer groups. Companies are frequently surprised to find that a small percentage of their products and customers contribute to a large percentage of their profits. This knowledge helps a company allocate its resources and effort.

Efficiency control involves micro-level analysis of the various elements of the marketing mix, including sales force, advertising, sales promotion, and distribution. For example, to understand its sales-force efficiency, a company may keep track of how many sales calls a repre-

sentative makes each day, how long each call lasts, and how much each call costs and generates in revenue. This type of analysis highlights areas in which companies can manage their marketing efforts in a more productive and cost-effective manner.

Strategic control processes allow managers to evaluate a company's marketing program from a critical long-term perspective. This involves a detailed and objective analysis of a company's organization and its ability to maximize its strengths and market opportunities. Companies can use two types of strategic control tools. The first, which a company uses to evaluate itself, is called a marketing-effectiveness rating review. In order to rate its own marketing effectiveness, a company examines its customer philosophy, the adequacy of its marketing information, and the efficiency of its marketing operations. It will also closely evaluate the strength of its marketing strategy and the integration of its marketing tactics.

The second evaluation tool is known as a marketing audit. This is a comprehensive, systematic, independent, and periodic analysis that a company uses to examine its strengths in relation to its current and potential market(s). Such an analysis is comprehensive because it covers all aspects of the marketing climate (unlike a functional audit, which analyzes one marketing activity), looking at both macro-environment factors (demographic, economic, ecological, technological, political, and cultural) and micro- or task-environment factors (markets, customers, competitors, distributors, dealers, suppliers, facilitators, and publics). The audit includes analyses of the company's marketing strategy, marketing organization, marketing systems, and marketing productivity. It must be systematic in order to provide concrete conclusions based on these analyses. To ensure objectivity, a marketing audit is best done by a person, department, or organization that is independent of the company or marketing program. Marketing audits should be done not only when the value of a company's current marketing plan is in question; they must be done periodically in order to isolate and solve problems before they arise.

(K.A.G./Jo.D.H./P.Ko.)

Customers

The elements that play a role in the marketing process can be divided into three groups: customers, distributors, and facilitators. In addition to interacting with one another, these groups must interact within a business environment that is affected by a variety of forces, including governmental, economic, and social influences.

In order to understand target customers, certain questions must be answered: Who constitutes the market segment? What do they buy and why? And how, when, and where do they buy? Knowing who constitutes the market segment is not simply a matter of knowing who uses a product. Often, individuals other than the user may participate in or influence a purchasing decision. Several individuals may play various roles in the decision-making process. For instance, in the decision to purchase an automobile for a small family business, the son may be the initiator, the daughter may be an influencer, the wife may be the decider, the purchasing manager may be the buyer, and the husband may be the user. In other words, the son may read in a magazine that businesses can save money and decrease tax liability by owning or leasing company transportation. He may therefore initiate the product search process by raising this issue at a weekly business meeting. However, the son may not be the best-qualified to gather and process information about automobiles, because the daughter worked for several years in the auto industry before joining the family business. Although the daughter's expertise and research efforts may influence the process, she may not be the key decision maker. The mother, by virtue of her position in the business and in the family, may make the final decision about which car to purchase. However, the family uncle may have good negotiation skills, and he may be the purchasing agent. Thus, he will go to different car dealerships in order to buy the chosen car at the best possible price. Finally, despite the involvement of all these individuals in the purchase process, none

Strategic
control
tools

of them may actually drive the car. It may be purchased so that the father may use it for his frequent sales calls. In other instances, an individual may handle more than one of these purchasing functions and may even be responsible for all of them. The key is that a marketer must recognize that different people have different influences on the purchase decision, and these factors must be taken into account in crafting a marketing strategy.

In addition to knowing to whom the marketing efforts are targeted, it is important to know which products target customers tend to purchase and why they do so. Customers do not purchase "things" as much as they purchase services or benefits to satisfy needs. For instance, a conventional oven allows users to cook and heat food. Microwave oven manufacturers recognized that this need could be fulfilled—and done so more quickly—with a technology other than conventional heating. By focusing on needs rather than on products, these companies were able to gain a significant share in the food cooking and heating market.

Knowledge of when, where, and how purchases are made is also useful. A furniture store whose target customers tend to make major purchases in the spring may send its mailings at the beginning of this season. A food vendor may set up a stand near the door of a busy office complex so that employees must pass the stand on their way to lunch. And a jeweler who knows that customers prefer to pay with credit cards may ensure that all major credit cards are accepted at the store. In other cases, marketers who understand specifics about buying habits and preferences also may try to alter them. Thus, a remotely situated wholesale store may use deeply discounted prices to lure customers away from the more conveniently located shopping malls.

Customers can be divided into two categories: consumer customers, who purchase goods and services for use by themselves and by those with whom they live; and business customers, who purchase goods and services for use by the organization for which they work. Although there are a number of similarities between the purchasing approaches of each type of customer, there are important differences as well.

CONSUMER CUSTOMERS

Factors influencing consumers. Four major types of factors influence consumer buying behaviour: cultural, social, personal, and psychological. Cultural factors have the broadest influence, because they constitute a stable set of values, perceptions, preferences, and behaviours that have been learned by the consumer throughout life. For example, in Western cultures consumption is often driven by a consumer's need to express individuality, while in Eastern cultures consumers are more interested in conforming to group norms. In addition to the influence of a dominant culture, consumers may also be influenced by several subcultures. In Quebec the dominant culture is French-speaking, but one influential subculture is English-speaking. Social class is also a subcultural factor: members of any given social class tend to share similar values, interests, and behaviours.

A consumer may interact with several individuals on a daily basis, and the influence of these people constitutes the social factors that impact the buying process. Social factors include reference groups—that is, the formal or informal social groups against which consumers compare themselves. Consumers may be influenced not only by their own membership groups but also by reference groups of which they wish to be a part. Thus, a consumer who wishes to be considered a successful white-collar professional may buy a particular kind of clothing because the people in this reference group tend to wear that style. Typically, the most influential reference group is the family. In this case, family includes the people who raised the consumer (the "family of orientation") as well as the consumer's spouse and children (the "family of procreation"). Within each group, a consumer will be expected to play a specific role or set of roles dictated by the norms of the group. Roles in each group generally are tied closely to status.

Personal factors include individual characteristics that, when taken in aggregate, distinguish the individual from others of the same social group and culture. These include age, life-cycle stage, occupation, economic circumstances, and lifestyle. A consumer's personality and self-conception will also influence his or her buying behaviour.

Finally, psychological factors are the ways in which human thinking and thought patterns influence buying decisions. Consumers are influenced, for example, by their motivation to fulfill a need. In addition, the ways in which an individual acquires and retains information will affect the buying process significantly. Consumers also make their decisions based on past experiences—both positive and negative.

Consumer buying tasks. A consumer's buying task is affected significantly by the level of purchase involvement. The level of involvement describes how important the decision is to the consumer; high involvement is usually associated with purchases that are expensive, infrequent, or risky. Buying also is affected by the degree of difference between brands in the product category. The buying task can be grouped into four categories based on whether involvement is high or low and whether brand differences are great or small.

Complex buying behaviour occurs when the consumer is highly involved with the purchase and when there are significant differences between brands. This behaviour can be associated with the purchase of a new home or of an advanced computer. Such tasks are complex because the risk is high (significant financial commitment), and the large differences among brands or products require gathering a substantial amount of information prior to purchase. Marketers who wish to influence this buying task must help the consumer process the information as readily as possible. This may include informing the consumer about the product category and its important attributes, providing detailed information about product benefits, and motivating sales personnel to influence final brand choice. For instance, realtors may offer consumers a book or a video featuring photographs and descriptions of each available home. And a computer salesperson is likely to spend time in the retail store providing information to customers who have questions.

Dissonance-reducing buying behaviour occurs when the consumer is highly involved but sees little difference between brands. This is likely to be the case with the purchase of a lawn mower or a diamond ring. After making a purchase under such circumstances, a consumer is likely to experience the dissonance that comes from noticing that other brands would have been just as good, if not slightly better, in some dimensions. A consumer in such a buying situation will seek information or ideas that justify the original purchase.

There are two types of low-involvement purchases. Habitual buying behaviour occurs when involvement is low and differences between brands are small. Consumers in this case usually do not form a strong attitude toward a brand but select it because it is familiar. In these markets, promotions tend to be simple and repetitive so that the consumer can, without much effort, learn the association between a brand and a product class. Marketers may also try to make their product more involving. For instance, toothpaste was at one time purchased primarily out of habit, but Procter and Gamble Co. introduced a brand, Crest toothpaste, that increased consumer involvement by raising awareness about the importance of good dental hygiene.

Variety-seeking buying behaviour occurs when the consumer is not involved with the purchase, yet there are significant brand differences. In this case, the cost of switching products is low, and so the consumer may, perhaps simply out of boredom, move from one brand to another. Such is often the case with frozen desserts, breakfast cereals, and soft drinks. Dominant firms in such a market situation will attempt to encourage habitual buying and will try to keep other brands from being considered by the consumer. These strategies reduce customer switching behaviour. Challenger firms, on the other hand, want consumers to switch from the market leader, so they

will offer promotions, free samples, and advertising that encourage consumers to try something new.

The consumer buying process. The purchase process is initiated when a consumer becomes aware of a need. This awareness may come from an internal source such as hunger or an external source such as marketing communications. Awareness of such a need motivates the consumer to search for information about options with which to fulfill the need. This information can come from personal sources, commercial sources, public or government sources, or the consumer's own experience. Once alternatives have been identified through these sources, consumers evaluate the options, paying particular attention to those attributes the consumer considers most important. Evaluation culminates with a purchase decision, but the buying process does not end here. In fact, marketers point out that a purchase represents the beginning, not the end, of a consumer's relationship with a company. After a purchase has been made, a satisfied consumer is more likely to purchase another company product and to say positive things about the company or its product to other potential purchasers. The opposite is true for dissatisfied consumers. Because of this fact, many companies continue to communicate with their customers after a purchase in an effort to influence post-purchase satisfaction and behaviour.

For example, a plumber may be motivated to consider buying a new set of tools because his old set of tools is getting rusty. To gather information about what kind of new tool set to buy, this plumber may examine the tools of a colleague who just bought a new set, read advertisements in plumbing trade magazines, and visit different stores to examine the sets available. The plumber then processes all the information collected, focusing perhaps on durability as one of the most important attributes. In making a particular purchase, the plumber initiates a relationship with a particular tool company. This company may try to enhance post-purchase loyalty and satisfaction by sending the plumber promotions about new tools.

BUSINESS CUSTOMERS

Business customers, also known as industrial customers, purchase products or services to use in the production of other products. Such industries include agriculture, manufacturing, construction, transportation, and communication, among others. They differ from consumer markets in several respects. Because the customers are organizations, the market tends to have fewer and larger buyers than consumer markets. This often results in closer buyer-seller relationships, because those who operate in a market must depend more significantly on one another for supply and revenue. Business customers also are more concentrated; for instance, in the United States more than half of the country's business buyers are concentrated in only seven states. Demand for business goods is derived demand, which means it is driven by a demand for consumer goods. Therefore, demand for business goods is more volatile, because variations in consumer demand can have a significant impact on business-goods demand. Business markets are also distinctive in that buyers are professional purchasers who are highly skilled in negotiating contracts and maximizing efficiency. In addition, several individuals within the business usually have direct or indirect influence on the purchasing process.

Factors influencing business customers. Although business customers are affected by the same cultural, social, personal, and psychological factors that influence consumer customers, the business arena imposes other factors that can be even more influential. First, there is the economic environment, which is characterized by such factors as primary demand, economic forecast, political and regulatory developments, and the type of competition in the market. In a highly competitive market such as airline travel, firms may be concerned about price and therefore make purchases with a focus on saving money. In markets where there is more differentiation among competitors—e.g., in the hotel industry—many firms may make purchases with a focus on quality rather than on price.

Second, there are organizational factors, which include the objectives, policies, procedures, structures, and systems

that characterize any particular company. Some companies are structured in such a way that purchases must pass through a complex system of checks and balances, while other companies allow purchasing managers to make more individual decisions. Interpersonal factors are more salient among business customers, because the participants in the buying process—perhaps representing several departments within a company—often have different interests, authority, and persuasiveness. Furthermore, the factors that affect an individual in the business buying process are related to the participant's role in the organization. These factors include job position, risk attitudes, and income.

The business buying process. The business buying process mirrors the consumer buying process, with a few notable exceptions. Business buying is not generally need-driven and is instead problem-driven. A business buying process is usually initiated when someone in the company sees a problem that needs to be solved or recognizes a way in which the company can increase profitability or efficiency. The ensuing process follows the same pattern as that of consumers, including information search, evaluation of alternatives, purchase decision, and post-purchase evaluation. However, in part because business purchase decisions require accountability and are often closely analyzed according to cost and efficiency, the process is more systematic than consumer buying and often involves significant documentation. Typically, a purchasing agent for a business buyer will generate documentation regarding product specifications, preferred supplier lists, requests for bids from suppliers, and performance reviews.

(K.A.G./Jo.D.H./P.Ko.)

Marketing intermediaries

Many producers do not sell products or services directly to consumers and instead use marketing intermediaries to execute an assortment of necessary functions to get the product to the final user. These intermediaries, such as middlemen (wholesalers, retailers, agents, and brokers), distributors, or financial intermediaries, typically enter into longer-term commitments with the producer and make up what is known as the marketing channel, or the channel of distribution. Manufacturers use raw materials to produce finished products, which in turn may be sent directly to the retailer, or, less often, to the consumer. However, as a general rule, finished goods flow from the manufacturer to one or more wholesalers before they reach the retailer and, finally, the consumer. Each party in the distribution channel usually acquires legal possession of goods during their physical transfer, but this is not always the case. For instance, in consignment selling, the producer retains full legal ownership even though the goods may be in the hands of the wholesaler or retailer—that is, until the merchandise reaches the final user or consumer.

Channels of distribution tend to be more direct—that is, shorter and simpler—in the less industrialized nations. There are notable exceptions, however. For instance, the Ghana Cocoa Marketing Board collects cacao beans in Ghana and licenses trading firms to process the commodity. Similar marketing processes are used in other West African nations. Because of the vast number of small-scale producers, these agents operate through middlemen who, in turn, enlist sub-buyers to find runners to transport the products from remote areas. Japan's marketing organization was, until the late 20th century, characterized by long and complex channels of distribution and a variety of wholesalers. It was possible for a product to pass through a minimum of five separate wholesalers before it reached a retailer.

Companies have a wide range of distribution channels available to them, and structuring the right channel may be one of the company's most critical marketing decisions. Businesses may sell products directly to the final customer, as Land's End, Inc., does with its mail-order goods and as is the case with most industrial capital goods. Or they may use one or more intermediaries to move their goods to the final user. The design and structure of consumer marketing channels and industrial marketing channels can be quite similar or vary widely.

Role of marketing after the purchase

Channel of distribution

The channel design is based on the level of service desired by the target consumer. There are five primary service components that facilitate the marketer's understanding of what, where, why, when, and how target customers buy certain products. The service variables are quantity or lot size (the number of units a customer purchases on any given purchase occasion), waiting time (the amount of time customers are willing to wait for receipt of goods), proximity or spatial convenience (accessibility of the product), product variety (the breadth of assortment of the product offering), and service backup (additional services such as delivery or installation provided by the channel). It is essential for the designer of the marketing channel—typically the manufacturer—to recognize the level of each service point that the target customer desires. A single manufacturer may service several target customer groups through separate channels, and therefore each set of service outputs for these groups could vary. One group of target customers may want elevated levels of service (that is, fast delivery, high product availability, large product assortment, and installation). Their demand for such increased service translates into higher costs for the channel and higher prices for customers. However, the prosperity of discount and warehouse stores demonstrates that customers are willing to accept lower service outputs if this leads to lower prices.

Channel functions and flows. In order to deliver the optimal level of service outputs to their target consumers, manufacturers are willing to allocate some of their tasks, or marketing flows, to intermediaries. As any marketing channel moves goods from producers to consumers, the marketing intermediaries perform, or participate in, a number of marketing flows, or activities. The typical marketing flows, listed in the usual sequence in which they arise, are collection and distribution of marketing research information (information), development and dissemination of persuasive communications (promotion), agreement on terms for transfer of ownership or possession (negotiation), intentions to buy (ordering), acquisition and allocation of funds (financing), assumption of risks (risk taking), storage and movement of product (physical possession), buyers paying sellers (payment), and transfer of ownership (title).

Each of these flows must be performed by a marketing intermediary for any channel to deliver the goods to the final consumer. Thus, each producer must decide who will perform which of these functions in order to deliver the service output levels that the target consumers desire. Producers delegate these flows for a variety of reasons. First, they may lack the financial resources to carry out the intermediary activities themselves. Second, many producers can earn a superior return on their capital by investing profits back into their core business rather than into the distribution of their products. Finally, intermediaries, or middlemen, offer superior efficiency in making goods and services widely available and accessible to final users. For instance, in overseas markets it may be difficult for an exporter to establish contact with end users, and various kinds of agents must therefore be employed. Because an intermediary typically focuses on only a small handful of specialized tasks within the marketing channel, each intermediary, through specialization, experience, or scale of operation, can offer a producer greater distribution benefits.

Management of channel systems. Although middlemen can offer greater distribution economy to producers, gaining cooperation from these middlemen can be problematic. Middlemen must continuously be motivated and stimulated to perform at the highest level. In order to gain such a high level of performance, manufacturers need some sort of leverage. Researchers have distinguished five bases of power: coercive (threats if the middlemen do not comply), reward (extra benefits for compliance), legitimate (power by position—rank or contract), expert (special knowledge), and referent (manufacturer is highly respected by the middlemen).

As new institutions emerge or products enter different life-cycle phases, distribution channels change and evolve. With these types of changes, no matter how well the chan-

nel is designed and managed, conflict is inevitable. Often this conflict develops because the interests of the independent businesses do not coincide. For example, franchisers, because they receive a percentage of sales, typically want their franchisees to maximize sales, while the franchisees want to maximize their profits, not sales. The conflict that arises may be vertical, horizontal, or multichannel in nature. When General Motors Corporation comes into conflict with its dealers, this is a vertical channel conflict. Horizontal channel conflict arises when a franchisee in a neighbouring town feels a fellow franchisee has infringed on its territory. Finally, multichannel conflict occurs when a manufacturer has established two or more channels that compete against each other in selling to the same market. For example, a major tire manufacturer may begin selling its tires through mass merchandisers, much to the dismay of its independent tire dealers.

Conflicting interests in the marketing channel

Merchandise distributors

WHOLESALESA

Wholesaling includes all activities required to sell goods or services to other firms, either for resale or for business use, usually in bulk quantities and at lower-than-retail prices. Wholesalers, also called distributors, are independent merchants operating any number of wholesale establishments. Wholesalers are typically classified into one of three groups: merchant wholesalers, brokers and agents, and manufacturers' and retailers' branches and offices.

Merchant wholesalers. Merchant wholesalers, also known as jobbers, distributors, or supply houses, are independently owned and operated organizations that acquire title ownership of the goods that they handle. There are two types of merchant wholesalers: full-service and limited-service.

Full-service wholesalers usually handle larger sales volumes; they may perform a broad range of services for their customers, such as stocking inventories, operating warehouses, supplying credit, employing salespeople to assist customers, and delivering goods to customers. General-line wholesalers carry a wide variety of merchandise, such as groceries; specialty wholesalers, on the other hand, deal with a narrow line of goods, such as coffee and tea, cigarettes, or seafood.

Limited-service wholesalers, who offer fewer services to their customers and suppliers, emerged in order to reduce the costs of service. There are several types of limited-service wholesalers. Cash-and-carry wholesalers usually handle a limited line of fast-moving merchandise, selling to smaller retailers on a cash-only basis and not delivering goods. Truck wholesalers or jobbers sell and deliver directly from their vehicles, often for cash. They carry a limited line of semiperishables such as milk, bread, and snack foods. Drop shippers do not carry inventory or handle the merchandise. Operating primarily in bulk industries such as lumber, coal, and heavy equipment, they take orders but have manufacturers ship merchandise directly to final consumers. Rack jobbers, who handle nonfood lines such as housewares or personal goods, primarily serve drug and grocery retailers. Rack jobbers typically perform such functions as delivery, shelving, inventory stacking, and financing. Producers' cooperatives—owned by their members, who are farmers—assemble farm produce to be sold in local markets and share profits at the end of the year.

In less-developed countries, wholesalers are often the sole or primary means of trade; they are the main elements in the distribution systems of many countries in Latin America, East Asia, and Africa. In such countries the business activities of wholesalers may expand to include manufacturing and retailing, or they may branch out into nondistributive ventures such as real estate, finance, or transportation. Until the late 1950s, Japan was dominated by wholesaling. Even relatively large manufacturers and retailers relied principally on wholesalers as their intermediaries. However, in the late 20th century, Japanese wholesalers have declined in importance. Even in the most highly industrialized nations, however, wholesalers remain essential to the operations of significant numbers of small retailers.

Brokers and agents. Manufacturers may use brokers and agents, who do not take title possession of the goods, in marketing their products. Brokers and agents typically perform only a few of the marketing flows, and their main function is to ease buying and selling—that is, to bring buyers and sellers together and negotiate between them. Brokers, most commonly found in the food, real estate, and insurance industries, may represent either a buyer or a seller and are paid by the party who hires them. Brokers often can represent several manufacturers of noncompeting products on a commission basis. They do not carry inventory or assume risk.

Manu-
facturers'
agents

Unlike merchant wholesalers, agent middlemen do not take legal ownership of the goods they sell; nor do they generally take physical possession of them. The three principal types of agent middlemen are manufacturers' agents, selling agents, and purchasing agents. Manufacturers' agents, who represent two or more manufacturers' complementary lines on a continuous basis, are usually compensated by commission. As a rule, they carry only part of a manufacturer's output, perhaps in areas where the manufacturer cannot maintain full-time salespeople. Many manufacturers' agents are businesses of only a few employees and are most commonly found in the furniture, electric, and apparel industries. Sales agents are given contractual authority to sell all of a manufacturer's output and generally have considerable autonomy to set prices, terms, and conditions of sale. Sometimes they perform the duties of a manufacturer's marketing department, although they work on a commission basis. Sales agents often provide market feedback and product information to the manufacturers and play an important role in product development. They are found in such product areas as chemicals, metals, and industrial machinery and equipment. Purchasing agents, who routinely have long-term relationships with buyers, typically receive, inspect, store, and ship goods to their buyers.

Manufacturers' and retailers' branches and offices. Wholesaling operations conducted by the sellers or buyers themselves rather than by independent wholesalers comprise the third major type of wholesaling. Manufacturers may engage in wholesaling through their sales branches and offices. This allows manufacturers to improve the inventory control, selling, and promotion flows. Numerous retailers also establish purchasing offices in major market centres such as Chicago and New York City that play a role similar to that of brokers and agents. The major difference is that they are part of the buyer's own organization.

RETAILERS

Retailing, the merchandising aspect of marketing, includes all activities required to sell directly to consumers for their personal, nonbusiness use. The firm that performs this consumer selling—whether it is a manufacturer, wholesaler, or retailer—is engaged in retailing. Retailing can take many forms: goods or services may be sold in person, by mail, telephone, television, or computer, or even through vending machines. These products can be sold on the street, in a store, or in the consumer's home. However, businesses that are classified as retailers secure the vast majority of their sales volume from store-based retailing.

The history of retailing. For centuries most merchandise was sold in marketplaces or by peddlers. In many countries, hawkers still sell their wares while traveling from one village to the next. Marketplaces are still the primary form of retail selling in these villages. This was also true in Europe until the Renaissance, when market stalls in certain localities became permanent and eventually grew into stores and business districts.

Retail chains are known to have existed in China several centuries before the Christian era and in some European cities in the 16th and 17th centuries. However, the birth of the modern chain store can be traced to 1859, with the inauguration of what is now the Great Atlantic & Pacific Tea Company, Inc. (A&P), in New York City. During the 15th and 16th centuries the Fugger family of Germany was the first to carry out mercantile operations of a chain-store variety. In 1670 the Hudson's Bay Company chartered its chain of outposts in Canada.

Department stores also were seen in Europe and Asia as early as the 17th century. The famous Bon Marché in Paris grew from a large specialty store into a full-fledged department store in the mid-1800s. By the middle of the 20th century, department stores existed in major U.S. cities, although small independent merchants still constitute the majority of retailers.

Shopping malls, a late 20th-century development in retail practices, were created to provide for a consumer's every need in a single, self-contained shopping area. Although they were first created for the convenience of suburban populations, they can now also be found on main city thoroughfares. A large branch of a well-known retail chain usually serves as a mall's retail flagship, which is the primary attraction for customers. In fact, few malls can be financed and built without a flagship establishment already in place.

Other mall proprietors have used recreation and entertainment to attract customers. Movie theatres, holiday displays, and live musical performances are often found in shopping malls. In Asian countries, malls also have been known to house swimming pools, arcades, and amusement parks. Hong Kong's City Plaza shopping mall includes one of the territory's two ice rinks. Some malls, such as the Mall of America in Bloomington, Minn., U.S., may offer exhibitions, sideshows, and other diversions.

Although there is a great variety of retail enterprises, with new types constantly emerging, they can be classified into three main types: store retailers, nonstore retailers, and retail organizations.

Store retailers. Several different types of stores participate in retail merchandising. The following is a brief description of the most important store retailers.

Specialty stores. A specialty store carries a deep assortment within a narrow line of goods. Furniture stores, florists, sporting-goods stores, and bookstores are all specialty stores. Stores such as Athlete's Foot (sports shoes only) and Tall Men (clothing for tall men) are considered superspecialty stores because they carry a very narrow product line.

Department stores. Department stores carry a wider variety of merchandise than most stores but offer these items in separate departments within the store. These departments usually include home furnishings and household goods, as well as clothing, which may be divided into departments according to gender and age. Department stores in western Europe and Asia also have large food departments, such as the renowned food court at Harrods in the United Kingdom. Departments within each store are usually operated as separate entities, each with its own buyers, promotions, and service personnel. Some departments, such as restaurants and beauty parlours, are leased to external providers.

Department stores generally account for less than 10 percent of a country's total retail sales, but they draw large numbers of customers in urban areas. The most influential of the department stores may even be trendsetters in various fields, such as fashion. Department stores such as Sears, Roebuck and Company have also spawned chain organizations. Others may do this through mergers or by opening branch units within a region or by expanding to other countries.

Supermarkets. Supermarkets are characterized by large facilities (15,000 to 25,000 square feet [1,394 to 2,323 square metres]) with more than 12,000 items, low profit margins (earning about 1 percent operating profit on sales), high volume, and operations that serve the consumer's total needs for items such as food (groceries, meats, produce, dairy products, baked goods) and household sundries. They are organized according to product departments and operate primarily on a self-service basis. Supermarkets also may sell wines and other alcoholic beverages (depending on local licensing laws) and clothing.

The first true supermarket was opened in the United States by Michael Cullin in 1930. His King Kullen chain of large-volume food stores was so successful that it encouraged the major food-store chains to convert their specialty stores into supermarkets. When compared with the conventional independent grocer, supermarkets generally

Importance
of depart-
ment stores

offered greater variety and convenience and often better prices as well. Consequently, in the two decades after World War II, the supermarket drove many small food retailers out of business, not only in the United States but throughout the world. In France, for example, the number of larger food stores grew from about 50 in 1960 to 4,700 in 1982, while the number of small food retailers fell from 130,000 to 60,000.

Convenience stores. Located primarily near residential areas, convenience stores are relatively small outlets that are open long hours and carry a limited line of high-turnover convenience products at high prices. Although many have added food services, consumers use them mainly for "fill-in" purchases, such as bread, milk, or miscellaneous goods.

Superstores. Superstores, hypermarkets, and combination stores are unique retail merchandisers. With facilities averaging 35,000 square feet, superstores meet many of the consumer's needs for food and nonfood items by housing a full-service grocery store as well as such services as dry cleaning, laundry, shoe repair, and cafeterias. Combination stores typically combine a grocery store and a drug store in one facility, utilizing approximately 55,000 square feet of selling space. Hypermarkets combine supermarket, discount, and warehousing retailing principles by going beyond routinely purchased goods to include furniture, clothing, appliances, and other items. Ranging in size from 80,000 to 220,000 square feet, hypermarkets display products in bulk quantities that require minimum handling by store personnel.

Discount stores. Selling merchandise below the manufacturer's list price is known as discounting. The discount store has become an increasingly popular means of retailing. Following World War II, a number of retail establishments in the United States began to pursue a high-volume, low-profit strategy designed to attract price-conscious consumers. A key strategy for keeping operating costs (and therefore prices) low was to locate in low-rent shopping districts and to offer minimal service assistance. This no-frills approach was used at first only with hard goods, or consumer durables, such as electrical household appliances, but it has since been shown to be successful with soft goods, such as clothing. This practice has been adopted for a wide variety of products, so that discount stores have essentially become department stores with reduced prices and fewer services. In the late 20th century, discount stores began to operate outlet malls. These groups of discount stores are usually located some distance away from major metropolitan areas and have facilities that make them indistinguishable from standard shopping malls. As they gained popularity, many discount stores improved their facilities and appearances, added new lines and services, and opened suburban branches. Coupled with attempts by traditional department stores to reduce prices in order to compete with discounters, the distinction between many department and discount stores has become blurred. Specialty discount operations have grown significantly in electronics, sporting goods, and books.

Off-price retailers. Off-price retailers offer a different approach to discount retailing. As discount houses tried to increase services and offerings in order to upgrade, off-price retailers invaded this low-price, high-volume sector. Off-price retailers purchase at below-wholesale prices and charge less than retail prices. This practice is quite different from that of ordinary discounters, who buy at the market wholesale price and simply accept lower margins by pricing their products below retail costs. Off-price retailers carry a constantly changing collection of overruns, irregulars, and leftover goods and have made their biggest forays in the clothing, footwear, and accessories industries. The three primary examples of off-price retailers are factory outlets, independent carriers, and warehouse clubs. Stocking manufacturers' surplus, discontinued, or irregular products, factory outlets are owned and operated by the manufacturer. Independent off-price retailers carry a rapidly changing collection of higher-quality merchandise and are typically owned and operated by entrepreneurs or divisions of larger retail companies. Warehouse (or wholesale) clubs operate out of enormous, low-cost facilities and

charge patrons an annual membership fee. They sell a limited selection of brand-name grocery items, appliances, clothing, and miscellaneous items at a deep discount. These warehouse stores, such as Wal-Mart-owned Sam's, Price Club, and Costco (in the United States), maintain low costs because they buy products at huge quantity discounts, use less labour in stocking, and typically do not make home deliveries or accept credit cards.

Nonstore retailers. Some retailers do not operate stores, and these nonstore businesses have grown much faster than store retailers. With some market observers predicting that by the year 2000 nonstore retailing will handle 30 percent of all general merchandise sold, nonstore channels may become a powerful force in the retailing industry. The major types of nonstore retailing are direct selling, direct marketing, and automatic vending.

Direct selling. This form of retailing originated several centuries ago and has mushroomed into a \$9 billion industry consisting of about 600 companies selling door-to-door, office-to-office, or at private-home sales meetings. The forerunners in the direct-selling industry include The Fuller Brush Company (brushes, brooms, etc.), Electrolux (vacuum cleaners), and Avon (cosmetics). In addition, Tupperware pioneered the home-sales approach, in which friends and neighbours gather in a home where Tupperware products are demonstrated and sold. Network marketing, a direct-selling approach similar to home sales, is also gaining prevalence in markets worldwide. Network marketing companies such as Amway and Shaklee reward their distributors not only for selling products but also for recruiting others to become distributors.

Direct marketing. Direct marketing is direct contact between a seller (manufacturer or retailer) and a consumer. Generally speaking, a seller can measure response to an offer because of its direct addressability. Although direct marketing gained wide popularity as a marketing strategy only in the late 20th century, it has been successfully utilized for more than a hundred years. The world's largest catalog houses—Sears, Roebuck and Company and Montgomery Ward & Co.—began as direct marketers in the late 1880s, selling their products solely by mail order. A century later, however, both companies were conducting most of their business in retail stores. Some department stores and specialty stores may supplement their store operations with direct-marketing transactions by mail or telephone. Mail-order firms grew rapidly in the 1950s and '60s in continental Europe, Great Britain, and certain other highly industrialized nations. Modern direct marketing is generally supported by advanced database technologies that track each customer's purchase behaviour. These technologies are used by established retail firms, such as Quelle and Neckermann in Germany, and are the foundation of mail-order businesses such as J. Crew, The Sharper Image, and L.L. Bean (all in the United States). Direct marketing is not a worldwide business phenomenon, however, because mail-order operations require infrastructure elements that are still lacking in many countries, such as efficient transportation networks and secure methods for transmitting payments.

Direct marketing has expanded from its early forms, among them direct mail and catalog mailings, to include such vehicles as telemarketing, direct-response radio and television, and electronic shopping. Unlike many other forms of promotion, a direct-marketing campaign is quantitatively measurable.

Automatic vending. Automatic vending is a unique area in nonstore merchandising because the variety of merchandise offered through automatic vending machines continues to grow. Initially, impulse goods with high convenience value such as cigarettes, soft drinks, candy, newspapers, and hot beverages were offered. However, a wide array of products such as hosiery, cosmetics, food snacks, postage stamps, paperback books, record albums, camera film, and even fishing worms are becoming available through machines.

Vending-machine operations are usually offered in sites owned by other businesses, institutions, and transportation agencies. They can be found in offices, gasoline stations, large retail stores, hotels, restaurants, and many other lo-

Growth of
nonstore
enterprise

cales. In Japan, vending machines now dispense frozen beef, fresh flowers, whiskey, jewelry, and even names of prospective dating partners. In Sweden, vending machines have developed as a supplementary channel to retail stores, where hours of business are restricted by law. High costs of manufacturing, installation, and operation have somewhat limited the expansion of vending-machine retailing. In addition, consumers typically pay a high premium for vended merchandise.

Retail organizations. While merchants can sell their wares through a store or nonstore retailing format, retail organizations can also structure themselves in several different ways. The major types of retail organizations are corporate chains, voluntary chains and retailer cooperatives, consumer cooperatives, franchise organizations, and merchandising conglomerates.

Corporate chains. Two or more outlets that have common ownership and control, centralized buying and merchandising operations, and similar lines of merchandise are considered corporate chain stores. Corporate chain stores appear to be strongest in the food, drug, shoe, variety, and women's clothing industries. Managed chain stores have a number of advantages over independently managed stores. Because managed chains buy large volumes of products, suppliers are willing to offer cost advantages that are not usually available to other stores. These savings can be passed on to consumers in the form of lower prices and better sales. In addition, because managed chains operate on such a large scale, they can hire more specialized and experienced personnel, who may be better able to take full advantage of purchasing and promotion opportunities. Chain stores also have the opportunity to take advantage of economies of scale in the areas of advertising, store design, and inventory control. However, a corporate chain may have disadvantages as well. Its size and bureaucracy often weaken staff members' personal interest, drive, creativity, and customer-service motivation.

Voluntary chains and retailer cooperatives. These are associations of independent retailers, unlike corporate chains. Wholesaler-sponsored voluntary chains of retailers who engage in bulk buying and collective merchandising are prevalent in many countries. True Value hardware stores represent this type of arrangement in the United States. In western Europe in the 1980s there were several large wholesaler-sponsored chains of retailers, each including more than 15,000 stores. These retail stores were located across 18 countries, each store using the same name and, as a rule, offering the same brands of products but remaining an independent enterprise. Wholesaler-sponsored chains offer the same types of services for their clients as do the financially integrated retail chains. Retailer cooperatives, such as ACE hardware stores, are grouped as independent retailers who establish a central buying organization and conduct joint promotion efforts.

Consumer cooperatives. Consumer cooperatives, or co-ops, are retail outlets that are owned and operated by consumers for their mutual benefit. The first consumer cooperative store was established in Rochdale, Eng., in 1844, and most co-ops are modeled after the same, original principles. They are based on open consumer membership, equal voting among members, limited customer services, and shared profits among members in the form of rebates generally related to the amounts of their purchases. Consumer cooperatives have gained widespread popularity throughout western and northern Europe, particularly in Denmark, Finland, Iceland, Norway, Sweden, and Great Britain. Co-ops typically emerge because community residents believe that local retailers' prices are too high or service is substandard.

Franchise organizations. Franchise arrangements are characterized by a contractual relationship between a franchiser (a manufacturer, wholesaler, or service organization) and franchisees (independent entrepreneurs who purchase the right to own and operate any number of units in the franchise systems). Typified by a unique product, service, business method, trade name, or patent, franchises have been prominent in many industries, including fast foods, video stores, health and fitness centres, hair salons, auto rentals, motels, and travel agencies. McDonald's Corpora-

tion is a prominent example of a franchise retail organization, with franchises all over the world.

Merchandising conglomerates. Merchandising conglomerates combine several diversified retailing lines and forms under central ownership, as well as integrate distribution and management of functions. Merchandising conglomerates are relatively free-form corporations. In the United States, Woolworth Corporation is considered a merchandising conglomerate because it operates Kinney shoe stores, Herald Square Stationers, Frame Scene, and Kids Mart.

(Jo.D.H./K.A.G./P.Ko.)

Marketing facilitators

Because marketing functions require significant expertise, it is often both efficient and effective for an organization to use the assistance of independent marketing facilitators. These are organizations and consultants whose sole or primary responsibility is to handle marketing functions. In many larger companies, all or some of these functions are performed internally. However, this is not necessary or justifiable in most companies, which usually require only part-time or periodic assistance from marketing facilitators. Also, most companies cannot afford to support the salaries and operating expenses required to maintain marketing facilitators as a permanent part of their staff. Furthermore, independent marketing contractors can be more effective than an internal department because nonemployee facilitators can have broader expertise and more objective perspectives. In addition, independent contractors often are more motivated to perform at high standards, because competition in the facilitator market is usually aggressive, and poor performance could mean lost business.

There are four major types of marketing facilitators: advertising agencies, market research firms, transportation firms, and warehousing firms.

ADVERTISING AGENCIES

Advertising agencies are responsible for initiating, managing, and implementing paid marketing communications. In addition, some agencies have diversified into other types of marketing communications, including public relations, sales promotion, interactive media, and direct marketing. Agencies typically consist of four departments: account management, a creative division, a research group, and a media planning department. Those in account management act as liaisons between the client and the agency, ensuring that client needs are communicated to the agency and that agency recommendations are clearly understood by the client. Account managers also manage the flow of work within the agency, making sure that projects proceed according to schedule. The creative department is where advertisements are conceived, developed, and produced. Artists, writers, and producers work together to craft a message that meets agency and client objectives. In this department, slogans, jingles, and logos are developed. The research department gathers and processes data about the target market and consumers. This information provides a foundation for the work of the creative department and account management. Media planning personnel specialize in selecting and placing advertisements in print and broadcast media.

Creating the advertising message

MARKET RESEARCH FIRMS

Market research firms gather and analyze data about customers, competitors, distributors, and other actors and forces in the marketplace. A large portion of the work performed by most market research firms is commissioned by specific companies for particular purposes. However, some firms also routinely collect a wide spectrum of data and then attempt to sell some or all of it to companies that may benefit from such information. For example, the A.C. Nielsen Co. in the United States specializes in supplying marketing data about consumer television viewing habits, and Information Resources, Inc. (IRI) has an extensive database regarding consumer supermarket purchases.

Marketing research may be quantitative, qualitative, or a combination of both. Quantitative research is numerically

Consumer-owned businesses

oriented, requires significant attention to the measurement of market phenomena, and often involves statistical analysis. For example, when a restaurant asks its customers to rate different aspects of its service on a scale from 1 (good) to 10 (poor), this provides quantitative information that may be analyzed statistically. Qualitative research focuses on descriptive words and symbols and usually involves observing consumers in a marketing setting or questioning them about their product or service consumption experiences. For example, a marketing researcher may stop a consumer who has purchased a particular type of detergent and ask him why that detergent was chosen. Qualitative and quantitative research each provides different insights into consumer behaviour, and research results are ordinarily more useful when the two methods are combined.

Market research can be thought of as the application of scientific method to the solution of marketing problems. It involves studying people as buyers, sellers, and consumers, examining their attitudes, preferences, habits, and purchasing power. Market research is also concerned with the channels of distribution, with promotion and pricing, and with the design of the products and services to be marketed.

TRANSPORTATION FIRMS

As a product moves from producer to consumer, it must often travel long distances. Many products consumed in the United States have been manufactured in another area of the world, such as Asia or Mexico. In addition, if the channel of distribution includes several firms, the product must be moved a number of times before it becomes accessible to consumers. A basic home appliance begins as a raw material (iron ore at a steel mill, for example) that is transported from a processing plant to a manufacturing facility.

Transportation firms assist marketers in moving products from one point in a channel to the next. An important matter of negotiation between companies working together in a channel is whether the sender or receiver of goods is responsible for transportation. Movement of products usually involves significant cost, risk, and time management. Thus, when firms consider a transportation option, they carefully weigh its dependability and price, frequency of operation, and accessibility. A firm that has its own transportation capabilities is known as a private carrier. There are also contract carriers, which are independent transportation firms that can be hired by companies on a long- or short-term basis. A common carrier provides services to any and all companies between predetermined points on a scheduled basis. The U.S. Postal Service is a common carrier, as are Federal Express and the Amtrak railway system.

WAREHOUSING FIRMS

Because products are not usually sold or shipped as soon as they are produced or delivered, firms require storage facilities. Two types of warehouses meet this need: storage warehouses hold goods for longer periods of time, and distribution warehouses serve as way stations for goods as they pass from one location to the next. Like the other marketing functions, warehouses can be wholly owned by firms, or space can be rented as needed. Although companies have more control over wholly owned facilities, warehouses of this sort can tie up capital and firm resources. Operations within warehouses usually require inspecting goods, tracking inventories, repackaging goods, shipping, and invoicing.

(K.A.G./P.Ko./Jo.D.H.)

Marketing in different sectors

Although the basic principles of marketing apply to all industries, the ways in which these principles are best applied can differ considerably based on the kind of product or service sold, the kind of buying behaviour associated with the purchase, and the sector (government, consumer goods, services, etc.).

The government market. This market consists of federal, state, and local governmental units that purchase or rent goods to fulfill their functions and responsibilities for

the public. Government agencies purchase a wide range of products and services, including helicopters, paintings, office furniture, clothing, alcohol, and fuel. Most of the agencies manage a significant portion of their own purchasing.

One prominent sector of the government market is the federal civilian buying establishment. In the United States this establishment consists of six categories: departments (e.g., the Department of Commerce), administration (e.g., the General Services Administration), agencies (e.g., the Federal Aviation Administration), boards (e.g., the Railroad Retirement Board), commissions (e.g., the Federal Communications Commission), and the executive office (e.g., the Office of Management and Budget). In addition there are several miscellaneous civilian buying establishments, such as, for example, the Tennessee Valley Authority.

Another governmental purchasing sector is the federal military buying establishment, represented in the United States by the Department of Defense, which purchases primarily through the Defense Supply Agency and the army, navy, and air force. The Defense Supply Agency operates six supply centres, which specialize in construction, electronics, fuel, personnel support, and industrial and general supplies.

Government purchasing procedures fall into two categories: the open bid and the negotiated contract. Under open-bid buying, the government disseminates very specific information about the products and services required and requests bids from suppliers. Contracts generally are awarded to the lowest bidder. In negotiated-contract buying, a government agency negotiates directly with one or more companies regarding a specific project or supply need. In most cases, contracts are negotiated for complex projects that involve major research-and-development costs and in matters where there is little effective competition.

Consumer-goods marketing. Consumer goods can be classified according to consumer shopping habits. Convenience goods are those that the customer purchases frequently, immediately, and with minimum effort. Tobacco products, soaps, and newspapers are all considered convenience goods, as are common staples like ketchup or pasta. Convenience-goods purchasing is usually based on habitual behaviour, where the consumer will routinely purchase a particular product. Some convenience goods, however, may be purchased impulsively, involving no habit, planning, or search effort. These goods, usually displayed near the cash register in a store in order to encourage quick choice and purchase, include candy, razors, and batteries. A slightly different type of convenience product is the emergency good, which is purchased when there is an urgent need. Such goods include umbrellas and snow shovels, and these are usually distributed at a wide variety of outlets so that they will be readily available when necessary.

A second type of product is the shopping good, which usually requires a more involved selection process than convenience goods. A consumer usually compares a variety of attributes, including suitability, quality, price, and style. Homogeneous shopping goods are those that are similar in quality but different enough in other attributes (such as price, brand image, or style) to justify a search process. These products might include automobile tires or a stereo or television system. Homogeneous shopping goods are often sold strongly on price.

With heterogeneous shopping goods, product features become more important to the consumer than price. Such is often the case with the purchase of major appliances, clothing, furniture, and high-tech equipment. In this situation, the item purchased must be a certain size or colour and must perform very specific functions that cannot be fulfilled by all items offered by every supplier. With goods of this sort, the seller has to carry a wide assortment to satisfy individual tastes and must have well-trained salespeople to provide both information and advice to consumers.

Specialty goods have particularly unique characteristics and brand identifications for which a significant group of buyers is willing to make a special purchasing effort. Examples include specific brands of fancy products, lux-

Civilian and military markets

ury cars, professional photographic equipment, and high-fashion clothing. For instance, consumers who favour merchandise produced by a certain shoe manufacturer or furniture maker will, if necessary, travel considerable distances in order to purchase that particular brand. In specialty-goods markets, sellers do not encourage comparisons between options; buyers invest time to reach dealers carrying the product desired, and these dealers therefore do not necessarily need to be conveniently located.

Finally, an unsought good is one that a consumer does not know about—or knows about but does not normally think of buying. New products, such as new frozen-food concepts or new communications equipment, are unsought until consumers learn about them through word-of-mouth influence or advertising. In addition, the need for unsought goods may not seem urgent to the consumer, and purchase is often deferred. This is frequently the case with life insurance, preventive car maintenance, and cemetery plots. Because of this, unsought goods require significant marketing efforts, and some of the more sophisticated selling techniques have been developed from the challenge to sell unsought goods.

Services marketing. A service is an act of labour or a performance that does not produce a tangible commodity and does not result in the customer's ownership of anything. Its production may or may not be tied to a physical product. Thus, there are pure services that involve no tangible product (as with psychotherapy), tangible goods with accompanying services (such as a computer software package with free software support), and hybrid product-services that consist of parts of each (for instance, restaurants are usually patronized for both their food and their service).

Services can be distinguished from products because they are intangible, inseparable from the production process, variable, and perishable. Services are intangible because they can often not be seen, tasted, felt, heard, or smelled before they are purchased. A person purchasing plastic surgery cannot see the results before the purchase, and a lawyer's client cannot anticipate the outcome of a case before the lawyer's work is presented in court. To reduce the uncertainty that results from this intangibility, marketers may strive to make their service tangible by emphasizing the place, people, equipment, communications, symbols, or price of the service. For example, consider the insurance slogans "You're in good hands with Allstate" or Prudential's "Get a piece of the Rock."

Services are inseparable from their production because they are typically produced and consumed simultaneously. This is not true of physical products, which are often consumed long after the product has been manufactured, inventoried, distributed, and placed in a retail store. Inseparability is especially evident in entertainment services or professional services. In many cases, inseparability limits the production of services because they are so directly tied to the individuals who perform them. This problem can be alleviated if a service provider learns to work faster or if the service expertise can be standardized and performed by a number of individuals (as H&R Block, Inc., has done with its network of trained tax consultants throughout the United States).

The variability of services comes from their significant human component. Not only do humans differ from one another, but their performance at any given time may differ from their performance at another time. The mechanics at a particular auto service garage, for example, may differ in terms of their knowledge and expertise, and each mechanic will have "good" days and "bad" days. Variability can be reduced by quality-control measures. These measures can include good selection and training of personnel and allowing customers to communicate dissatisfaction (e.g., through customer suggestion and complaint systems) so that poor service can be detected and corrected.

Finally, services are perishable because they cannot be stored. Because of this, it is difficult for service providers to manage anything other than steady demand. When demand increases dramatically, service organizations face the problem of producing enough output to meet customer needs. When a large tour bus unexpectedly arrives

at a restaurant, its staff must rush to meet the demand, because the food services (taking orders, making food, taking money, etc.) cannot be "warehoused" for such an occasion. To manage such instances, companies may hire part-time employees, develop efficiency routines for peak demand occasions, or ask consumers to participate in the service-delivery process. On the other hand, when demand drops off precipitously, service organizations are often burdened with a staff of service providers who are not performing. Organizations can maintain steady demand by offering differential pricing during off-peak times, anticipating off-peak hours by requiring reservations, and giving employees more flexible work shifts.

(K.A.G./Jo.D.H./P.Ko.)

Business marketing. Business marketing, sometimes called business-to-business marketing or industrial marketing, involves those marketing activities and functions that are targeted toward organizational customers. This type of marketing involves selling goods (and services) to organizations (public and private) to be used directly or indirectly in their own production or service-delivery operations. Some of the major industries that comprise the business market are construction, manufacturing, mining, transportation, public utilities, communications, and distribution. One of the key points that differentiates business from consumer marketing is the magnitude of the transactions. For example, in the mid-1990s, a Boeing 747 airliner, selling for about \$155 million, could take up to four years to manufacture and deliver once the order was placed. Often, a major airline company will order several aircraft at one time, making the purchase price as high as a billion dollars.

Customers for industrial goods can be divided into three groups: user customers, original-equipment manufacturers, and resellers. User customers make use of the goods they purchase in their own businesses. An automobile manufacturer, for example, might purchase a metal-stamping press to produce parts for its vehicles. Original-equipment manufacturers incorporate the purchased goods into their final products, which are then sold to final consumers (e.g., the manufacturer of television receivers buys tubes and transistors). Industrial resellers are middlemen—essentially wholesalers but in some cases retailers—who distribute goods to user customers, to original-equipment manufacturers, and to other middlemen. Industrial-goods wholesalers include mill-supply houses, steel warehouses, machine-tool dealers, paper jobbers, and chemical distributors.

Nonprofit marketing. Marketing scholars began exploring the application of marketing to nonprofit organizations in 1969. Since then, nonprofit organizations have increasingly turned to marketing for growth, funding, and prosperity.

Although it is difficult to define "nonprofit" organizations because of the existence of a number of quasi-governmental organizations, a study in the mid-1990s found more than one million private, nonprofit organizations in the United States. Some experts believe that the way to distinguish between organizations is according to their sources of funding. The three major sources are profits, government revenues (such as grants or taxes), and voluntary donations. In addition, a legally defined nonprofit organization is one that has been granted tax-exempt status by the Internal Revenue Service. However, while nonprofit groups can be defined legally, it is more helpful to focus on the specific marketing activities that need to be performed within the organization's environment. Museums, hospitals, universities, and churches are all examples of nonprofit organizations. Although many individuals may believe that nonprofit organizations have only a small impact on the economy, the operating expenditures of private nonprofit organizations now represent a significant percentage of the U.S. gross national product. In addition, many of these are substantial enterprises. For example, Girl Scout cookies, sold by Girl Scouts of America, constitute 10 percent of all cookies sold in the United States.

Social marketing. Social marketing employs marketing principles and techniques to advance a social cause, idea,

Increasing
service
production

or behaviour. It entails the design, implementation, and control of programs aimed at increasing the acceptability of a social idea or practice that would benefit the adoptors or society. Social ideas can take the form of beliefs, attitudes, and values, such as human rights. Whether social marketers are promoting ideas or social practices, their ultimate goal is to alter behaviour. In order to accomplish this behaviour change, social marketers set measurable objectives, research their target group's needs, target their "products" to these particular "consumers," and effectively communicate their benefits. In addition, social-marketing organizations have to be constantly aware of changes in their environments and must be able to adapt to these changes.

Place marketing. Place marketing employs marketing principles and techniques to advance the appeal and viability of a place (town, city, state, region, or nation) to tourists, businesses, investors, and residents. Among the "place sellers" are economic development agencies, tourist promotion agencies, and mayors' offices. Place sellers must gain a deep understanding of how place buyers make their purchasing decisions. Place-marketing activities can be found in both the private and public sectors at the local, regional, national, and international levels. They can range from activities involving downtrodden cities trying to attract businesses to vacation spots seeking to attract tourists. In implementing these marketing activities, each locale must adapt to external shocks and forces beyond its control (intergovernmental power shifts, increasing global competition, and rapid technological change) as well as to internal forces and decline cycles. (Jo.D.H./P.Ko.)

Economic and social aspects of marketing

Sometimes criticized for its impact on personal economic and social well-being, marketing has been said to affect not only individual consumers but also society as a whole. This section briefly examines some of the criticisms raised and how governments, individuals, and marketers have addressed them.

Marketing and individual welfare. Criticisms have been leveled against marketers, claiming that some of their practices may damage individual welfare. While this may be true in certain circumstances, it is important to recognize that, if a business damages individual welfare, it cannot hope to continue in the marketplace for long. As a consequence, most unfavourable views of marketing are criticisms of poor marketing, not of strategically sound marketing practices.

Others have raised concerns about marketing by saying that it increases prices by encouraging excessive markups. Marketers recognize that consumers may be willing to pay more for a product—such as a necklace from Tiffany and Co.—simply because of the associated prestige. This not only results in greater costs for promotion and distribution, but it allows marketers to earn profit margins that may be significantly higher than industry norms. Marketers counter these concerns by pointing out that products provide not only functional benefits but symbolic ones as well. By creating a symbol of prestige and luxury, Tiffany's offers a symbolic benefit that, according to some consumers, justifies the price. In addition, brands may symbolize not only prestige but also quality and functionality, which gives consumers greater confidence when they purchase a branded product. Finally, advertising and promotions are often very cost-effective methods of informing the general public about items and services that are available in the marketplace.

A few marketers have been accused of using deceptive practices, such as misleading promotional activities or high-pressure selling. These deceptive practices have given rise to legislative and administrative remedies, including guidelines offered by the Federal Trade Commission (FTC) regarding advertising practices, automatic 30-day guarantee policies by some manufacturers, and "cooling off" periods during which a consumer may cancel any contract signed. In addition, professional marketing associations, such as the Direct Marketing Association, have promulgated a set of professional standards for their industry.

Marketing and societal welfare. Concern also has been raised that some marketing practices may encourage excessive interest in material possessions, create "false wants," or promote the purchase of nonessential goods. For example, in the United States, children's Saturday morning television programming came under fire for promoting materialistic values. The Federal Communications Commission (FCC) responded in the early 1990s by regulating the amount of commercial time per hour. In many of these cases, however, the criticisms overstate the power of marketing communications to influence individuals and portray members of the public as individuals unable to distinguish between a good decision and a bad one. In addition, such charges cast marketing as a cause of social problems when often the problems have much deeper societal roots.

Marketing activity also has been sometimes criticized because of its control by strong private interests and its neglect of social and public concern. While companies in the cigarette, oil, and alcohol industries may have significant influence on legislation, media, and individual behaviour, organizations that focus on environmental, health, or education concerns are not able to wield such influence and often fail to receive appropriate recognition for their efforts. While there is clearly an imbalance of power between private interests and public ones, in the late 20th century, private companies have received more praise for their marketing efforts for social causes.

(P.Ko./K.A.G./Jo.D.H.)

Marketing's contributions to individuals and society. Although some have questioned the appropriateness of the marketing philosophy in an age of environmental deterioration, resource shortages, world hunger and poverty, and neglected social services, numerous firms are commendably satisfying individual consumer demands as well as acting in the long-term interests of the consumer and society. These dual objectives of many of today's companies have led to a broadening of the "marketing concept" to become the "societal marketing concept." Generating customer satisfaction while at the same time attending to consumer and societal well-being in the long run are the core concepts of societal marketing.

In practicing societal marketing, marketers try to balance company profits, consumer satisfaction, and public interest in their marketing policies. Many companies have achieved success in adopting societal marketing. Two prominent examples are The Body Shop International PLC, based in England, and Ben & Jerry's Homemade Inc., which produces ice cream and is based in Vermont. Body Shop's cosmetics and personal hygiene products, based on natural ingredients, are sold in recycled packaging. The products are formulated without animal testing, and a percentage of profits each year is donated to animal rights groups, homeless shelters, Amnesty International, rain-forest preservation groups, and other social causes. Ben & Jerry's donates a percentage of its profits to help alleviate social and environmental problems. The company's corporate concept focuses on "caring capitalism," which involves the product as well as social and economic missions.

Marketing has had many other positive benefits for individuals and society. It has helped accelerate economic development and create new jobs. It has also contributed to technological progress and enhanced consumers' choices.

(P.Ko./Jo.D.H.)

BIBLIOGRAPHY. The most widely used textbook is PHILIP KOTLER, *Marketing Management: Analysis, Planning, Implementation, and Control*, 8th ed. (1994). ROBERT BARTELS, *History of Marketing Thought*, 3rd ed. (1988), provides an overview of marketing through the years and contains an extensive bibliography. A more conceptual and theoretical treatment may be found in the work by WROE ALDERSON, *Dynamic Marketing Behavior: A Functional Theory of Marketing* (1965).

Various strategic and tactical aspects of marketing are explored in the following studies: STEVEN P. SCHNAARS, *Marketing Strategy: A Customer-Driven Approach* (1991); JACK TROUT and AL REIS, *Positioning: The Battle for Your Mind*, rev. ed. (1986); THOMAS T. NAGLE and REED K. HOLDEN, *The Strategy and Tactics of Pricing*, 2nd ed. (1994); GILBERT A. CHURCHILL, JR., NEIL M. FORD, and ORVILLE C. WALKER, JR., *Sales Force*

Management, 4th ed. (1993); DON E. SCHULTZ, STANLEY I. TANNENBAUM, and ROBERT F. LAUTERBORN, *Integrated Marketing Communications* (1992); MARY LOU ROBERTS and PAUL D. BERGER, *Direct Marketing Management* (1989); DAVID A. AAKER, *Strategic Market Management*, 3rd ed. (1991); GLEN L. URBAN and JOHN HAUSER, *Design and Marketing of New Products*, 2nd ed. (1993); and STAN RAPP and THOMAS L. COLLINS, *Beyond Maximarketing* (1994). The importance of marketing intermediaries is outlined by LOUIS W. STERN and ADEL I. EL-ANSARY, *Marketing Channels*, 4th ed. (1992).

THEODORE LEVITT, *The Marketing Imagination*, new, expanded ed. (1986); and REGIS MCKENNA, *Relationship Marketing: Successful Strategies for the Age of the Customer* (1991), discuss the evolving discipline of marketing.

The role of marketing research is investigated in the extended work by GILBERT A. CHURCHILL, JR., *Marketing Research: Methodological Foundations*, 5th ed. (1991). The role the consumer plays in the marketing process is examined in MICHAEL R. SOLOMON, *Consumer Behavior*, 2nd ed. (1994); and DAVID A. AAKER and GEORGE S. DAY, *Consumerism: Search for Consumer Interest*, 4th ed. (1982).

Advertising's role in the marketing process is explored in DAVID OGILVY, *Ogilvy on Advertising* (1983); JOHN LYONS, *Guts: Advertising from the Inside Out* (1987); and WILLIAM WELLS, JOHN BURNETT, and SANDRA MORIARTY, *Advertising: Principles and Practice*, 2nd ed. (1992).

Marketing in several different sectors is dealt with in the following selected works: DAVID A. AAKER and ALEXANDER L.

BIEL, *Brand Equity & Advertising* (1993); DAVID A. AAKER, *Managing Brand Equity* (1991); CHRISTOPHER H. LOVELOCK, *Services Marketing*, 2nd ed. (1991), and *Managing Services: Marketing, Operations, and Human Resources*, 2nd ed. (1991); MICHAEL D. HUTT and THOMAS W. SPEH, *Business Marketing Management: A Strategic View of Industrial and Organizational Markets*, 4th ed. (1992); PHILIP KOTLER and ALAN R. ANDREASEN, *Strategic Marketing for Nonprofit Organizations*, 4th ed. (1991); PHILIP KOTLER and ROBERTA N. CLARKE, *Marketing for Healthcare Organizations* (1986); PHILIP KOTLER and EDUARDO ROBERTO, *Social Marketing: Strategies for Changing Public Behavior* (1989); PHILIP KOTLER, DONALD H. HEIDER, and IRVING REIN, *Marketing Places: Attracting Investment, Industry, and Tourism to Cities, States, and Nations* (1993); and MICHAEL R. CZINKOTA and ILKKA A. RONKAINEN, *International Marketing*, 3rd ed. (1993).

Current marketing trends are reported in a number of trade journals and newspapers, including *Marketing News* (biweekly); *Marketing Management* (quarterly); *Journal of Retailing* (quarterly); *Advertising Age* (weekly); *Business Marketing* (monthly); *Harvard Business Review* (bimonthly); *Wall Street Journal* (daily); *Business Week* (weekly); *Fortune* (biweekly); *California Management Review* (quarterly); and *Business Horizons* (bimonthly). More scholarly research journals include *Journal of Marketing Research* (quarterly); *Journal of Marketing* (quarterly); *Journal of Consumer Research* (quarterly); and *Journal of the Academy of Marketing Science* (quarterly).

(Jo.D.H./P.Ko./K.A.G.)

Markets

Markets in the most literal and immediate sense are places in which things are bought and sold. In the modern industrial system, however, the market is not a place; it has expanded to include the whole geographical area in which sellers compete with each other for customers. Alfred Marshall, whose *Principles of Economics* (first published in 1890) was for long an authority for English-speaking economists, based his definition of the market on that of the French economist A. Cournot:

Economists understand by the term Market, not any particular market place in which things are bought and sold, but the whole of any region in which buyers and sellers are in such free intercourse with one another that the prices of the same goods tend to equality easily and quickly.

To this Marshall added:

The more nearly perfect a market is, the stronger is the tendency for the same price to be paid for the same thing at the same time in all parts of the market.

The concept of the market as defined above has to do primarily with more or less standardized commodities, for example, wool or automobiles. The word market is also used in contexts such as the market for real estate or for old masters; and there is the "labour market," although a contract to work for a certain wage differs from a sale of

goods. There is a connecting idea in all of these various usages—namely, the interplay of supply and demand.

Most markets consist of groups of intermediaries between the first seller of a commodity and the final buyer. There are all kinds of intermediaries, from the brokers in the great produce exchanges down to the village grocer. They may be mere dealers with no equipment but a telephone, or they may provide storage and perform important services of grading, packaging, and so on. In general, the function of a market is to collect products from scattered sources and channel them to scattered outlets. From the point of view of the seller, dealers channel the demand for his product; from the point of view of the buyer, they bring supplies within his reach.

There are two main types of markets for products, in which the forces of supply and demand operate quite differently, with some overlapping and borderline cases. In the first, the producer offers his goods and takes whatever price they will command; in the second, the producer sets his price and sells as much as the market will take. In addition, along with the growth of trade in goods, there has been a proliferation of financial markets, including securities exchanges and money markets.

The article is divided into the following sections:

The market in economic doctrine and history	509
Market theory	509
The abstract nature of traditional market theory	
Modifications of the theory	
The historical development of markets	510
The origin of markets	
Markets under Socialism	
Markets for primary products	511
Commodity markets	511
Futures markets	511
Economic functions of the futures contract	
The theory and practice of hedging	
Important futures markets	
Markets for manufactures	513
Financial markets	513
Securities trading	514

Types of corporate securities	
The marketing of new issues	
The machinery of securities trading	
The structure of demand for securities	
Dynamics of the securities markets	
Money market	519
Banks and the money market	
The international money market	
The U.S. money market	
The British money market	
The money markets of other countries	
The markets and social welfare	523
The politics of the market	523
Market Socialism	524
Bibliography	524

The market in economic doctrine and history

MARKET THEORY

The abstract nature of traditional market theory. The key to the modern concept of the market may be found in the famous observation of the 18th-century British economist Adam Smith that "The division of labour depends upon the extent of the market." He foresaw that modern industry depended for its development upon an extensive market for its products. The factory system developed out of trade in cotton textiles, when merchants, discovering an apparently insatiable worldwide market, became interested in increasing production in order to have more to sell. The factory system led to the use of power to supplement human muscle, followed in turn by the application of science to technology, which in an ever-accelerating spiral has produced the scope and complexity of modern industry.

The economic theory of the late 19th century, which is still influential in academic teaching, was, however, concerned with the allocation of existing resources between different uses rather than with technical progress. This theory was highly abstract. The concept of the market was most systematically worked out in a general equilibrium system developed by the French economist Léon Walras, who was strongly influenced by the theoretical physics of his time. His system of mathematical equations was ingenious, but there are two serious limitations to the mechan-

ical analogy upon which they were based: it omitted the factor of time—the effect upon peoples' present behaviour of their expectations about the future; and it ignored the consequences for the human beings concerned of the distribution of purchasing power among them. Though economists have always admitted the abstract nature of the theory, they generally have accepted the doctrine that the free play of market forces tended to bring about full employment and an optimum allocation of resources. On this view, unemployment could only be caused by wages being too high. This doctrine was still influential in the Great Depression of the 1930s.

Modifications of the theory. The change in view that was to become known as the Keynesian Revolution was largely an escape to common sense, as opposed to abstract theory. In a private-enterprise economy, investment in industrial installations and housing construction is aimed at profitability in the future. Because investment therefore depends upon expectations, unfavourable expectations tend to fulfill themselves—when investment outlay falls off, workers become unemployed; incomes fall, purchases fall, unemployment spreads to the consumer goods industries, and receipts are reduced all the more. The operation of the market thus generates instability. The market may also generate instability in an upward direction. A high level of effective demand leads to a scarcity of labour; rising wages raise both costs of production and incomes so that there is a general tendency to inflation.

The turn
toward
realism

While the English economist John Maynard Keynes was attacking the concept of equilibrium in the market as a whole, the notion of equilibrium in the market for particular commodities was also being undermined. Traditional theory had conceived of a group of producers as operating in a perfect market for a single commodity; each produced only a small part of the whole supply; for each, the price was determined by the market; and each maximized its profits by selling only as much as would make marginal cost equal to price—that is to say, only so much that it would to proceed. Each firm worked its plant up to capacity—*i.e.*, to the point where profitability was limited by rising costs. This state of affairs, known as “perfect competition,” is quite contrary to the general run of business experience, particularly in bad times when under-capacity working is prevalent. A theory of imperfect competition was invented to reconcile the traditional theory with under-capacity working but was attacked as unrealistic. The upshot was a general recognition that strict profit maximizing is impossible in conditions of uncertainty; that prices of manufactures are generally formed by adding a margin to direct costs, large enough to yield a profit at less than capacity sales; and that an increase in capacity generally has to be accompanied by a selling campaign to ensure that it will be used at a remunerative level.

Once it is recognized that competition is never perfect in reality, it becomes obvious that there is great scope for individual variations in the price policy of firms. No precise generalization is possible. The field is open for study of what actually happens, and exploration is going on. Meanwhile, however, textbook teaching often continues to seek refuge in the illusory simplicity of the traditional theory of market behaviour.

THE HISTORICAL DEVELOPMENT OF MARKETS

Economies without markets

History and anthropology provide many examples of economies based neither on markets nor on commerce. An exchange of gifts between communities with different resources, for example, may resemble trade, particularly in diversifying consumption and encouraging specialization in production, but subjectively it has a different meaning. Honour lies in giving; receiving imposes a burden. There is competition to see who can show the most generosity, not who can make the biggest gain. Another kind of non-commercial exchange was the payment of tribute, or dues, to a political authority, which then distributed what it had collected. On this basis, great, complex, and wealthy civilizations have arisen in which commerce was almost entirely unknown: the network of supply and distribution was operated through the administrative system. Herodotus remarked that the Persians had no marketplaces.

The distinguishing characteristic of commerce is that goods are offered not as a duty or for prestige or out of neighbourly kindness but in order to acquire purchasing power. It is clearly a convenience to all parties to have a single generally established currency-commodity. Once a commodity is acceptable as money, its use to store purchasing power overshadows its use for its original purpose; it ceases to be a commodity like any other and becomes the very embodiment of value.

The origin of markets. Markets as centres of commerce seem to have had three separate points of origin. The first was in rural fairs. A typical cultivator fed his family and paid the landlord and the moneylender from his chief crop. He had sidelines that provided salable products, and he had needs that he could not satisfy at home. It was then convenient for him to go to a market where many could meet to sell and buy.

The second point was in service to the landlords. Rent, essentially, was paid in grain; even when it was translated into money, sales of grain were necessary to supply the cultivator with funds to meet his dues. Payment of rent was a one-way transaction, imposed by the landlord. In turn, the landlord used the rents to maintain his warriors, clients, and artisans, and this led to the growth of towns as centres of trade and production. An urban class developed with a standard of life enabling its members to cater to each other as well as to the landlords and officials.

The third, and most influential, origin of markets was in international trade. From early times, merchant adventurers (the Phoenicians, the Arabs) risked their lives and their capital in carrying the products of one region to another. The importance of international trade for the development of the market system was precisely that it was carried on by third parties. Within a settled country, commercial dealings were restrained by considerations of rights, obligations, and proper behaviour. In medieval Europe, for example, dealings were regulated in the main by the concept of the “just price,” that is, a system of valuations that assured the producers and merchants an income sufficient to maintain life at a level suited to their respective positions in society. But in trade in which the dealer is not subject to any obligation at either end, no holds are barred; purely commercial principles have free play. It was in trade (for instance, the export of English wool to the weavers of Italy) that the commercial principle undermined feudal conceptions of rights and duties. As Adam Smith observed, a great leap occurred when trade released the forces of industrial production.

Throughout history the relations between the trader and the producer have changed with the development of technique and with changes in the economic power of the parties. The 19th century was the heyday of the import-export merchant. Traders from a metropolitan country could establish themselves in a foreign centre, become experts on its needs and possibilities, and deal with a great variety of producers and customers, on a relatively small scale with each. With the growth of giant corporations, the scope of the merchant narrowed; his functions were largely taken over by the sales departments of the industrial concerns. Nowadays it is common to hold international fairs at which industrial products are displayed for inspection by customers, a grand and glorified version of the village market; the business, however, consists in placing orders rather than buying on the spot and carrying merchandise home. The function of the independent wholesaler, like that of the merchant, has declined as great retail businesses have grown to a scale whereby they can deal directly with manufacturers; but specialized exchanges for primary commodities are still important.

Markets under Socialism. Markets are essential to the free enterprise system; they grew and spread along with it. The propensity “to truck, barter, and exchange one thing for another” (in Adam Smith’s words) was exalted into a principle of civilization by the doctrine of *laissez-faire*, which taught that the pursuit of self-interests by the individual would be to the benefit of society as a whole. In the Soviet Union and other Socialist countries, a different kind of economy existed and a different ideology was dominant. There were two interlocking systems in the economy of the Soviet Union: one for industry and one for agriculture; and the same pattern was followed, with variations, in the other Socialist countries. Industrially, all equipment and materials were owned by the state, and production was directed according to a central plan. In theory, payments to workers were thought of as their share of the total production of the economy; in practice, however, the system of wages was very much like that in capitalist industry except that rates as a rule were set by decree and the managers of enterprises had little scope for bargaining. Workers might move around looking for jobs, but there was no “labour market” in the capitalist sense. Materials and equipment were distributed among enterprises by the state planning offices. (Faulty planning gave rise to intermediaries who operated between enterprises, but this is not at all the same thing as the highly developed markets in materials, components, and equipment that exist under capitalism.)

Consumption goods, on the other hand, were distributed to Soviet households through a retail market. Though some Socialist idealists, regarding buying and selling as the essence of capitalism, have advocated that money should be abolished altogether, in a large community it has proved to be most convenient to provide incomes in the form of generalized purchasing power and to allow each to choose what he pleases from whatever goods are available. Classical economists usually assert that the advantage of the

International markets

Markets in Communist-governed countries

retail market system is that it runs itself without excessive regulation; consumers who go shopping are in charge of their own money and need account to no one for what they do with it. Retail markets in the Soviet economy differed from those in capitalist economies in that, while in both systems the buyer is in this sense a principal, the seller in the Soviet model was an agent. Retailers and manufacturers all served as agents of the same authority—the central plan. Rather than making it their business to woo and cajole the customer, sellers threw supplies into the shops in a somewhat arbitrary way and customers would search for what they wanted.

Soviet agriculture was organized on principles quite different from those operative for manufacturing. Collective farms, though managed in an authoritarian way, were like cooperatives in which members shared in the income of their farm in respect to the “work points” each could earn. The value of a work point was affected by the prices set for the products of the farm, and these were politically, rather than only economically, determined. In the Western industrial economies, there is also a political element involved in the setting of agricultural prices; generally the problem here is to prevent excess production from driving prices too low. For the Soviets, the problem was the opposite. There, agricultural output failed to expand rapidly enough to keep pace with the requirements of the growing industrial labour force, and prices were therefore kept down so that they would not be unfavourable to the industrial sector. At the same time, individual members of the collective farms were permitted to sell the produce of their household plots on a free market. In this specific market, the peasant was as much a principal as the buyer.

In China, cooperative farms established after 1949 were much more genuinely cooperatives than were those in the Soviet Union, and trade with the cities in China is organized through a kind of Socialist wholesaling. City authorities place contracts with neighbouring farms, specifying prices, varieties, quantities, and delivery dates, and then direct the supplies to retail outlets, which are part of the Socialist economy. A similar system controls trade in manufactured consumer goods. Through the retail shops, the authorities monitor demand and guide supply as far as possible to meet it by the contracts that they place with the Socialist manufacturers. By adapting the wholesale trade to its own requirements, the Chinese economy seems to have avoided some of the difficulties that the Soviets encountered.

An example of socialism without a formal market was seen in the early days of the cooperative settlements known as *kibbutzim* in Israel, where cultivators shared the proceeds of their work without any distinction of individual incomes. (Because a *kibbutz* could trade with the surrounding market economy, its members were not confined to consuming only the produce of their own soil.) At the outset some of the *kibbutzim* carried the objection to private property so far that a man who gave a shirt to the laundry received back just some other shirt. But to dispense altogether with market relationships is apparently possible only in a small community in which all share a common ideal, and the austere standards of the original *kibbutzim* have softened somewhat with growing prosperity; but they still maintain a small-scale example of economic efficiency without commercial incentives.

Markets for primary products

COMMODITY MARKETS

The general run of agricultural commodities is produced under competitive conditions by relatively small-scale cultivators scattered over a large area. The final purchasers are also scattered, and centres of consumption are distant from regions of production. The dealer, therefore, since he is indispensable, is in a stronger economic position than the seller. This situation is markedly true when the producer is a peasant who lacks both commercial knowledge and finance so that he is obliged to sell as soon as his harvest comes in; it is true also, though to a lesser extent, of the capitalist plantation for which the only source of earnings is a particular specialized product. In this kind of

business, both demand and supply are said to be inelastic in the short run—that is, a fall in price does not have much effect in increasing purchases and a rise in price cannot quickly increase supplies. Supplies are subject to natural variations, weather conditions, pests, and so forth; and demand varies with the level of activity in the centres of industry and with changes in tastes and technical requirements. Under a regime of unregulated competition such markets are, therefore, tormented with continual fluctuations in prices and volume of business. Though dealers may mitigate this to some extent by building up stocks when prices are low and releasing them when demand is high, such buying and selling often turns into speculation, which tends to exacerbate the fluctuations.

The behaviour of primary commodity markets is a serious matter when whole communities depend upon a single commodity for income or for employment and wages. The agricultural communities that form part of an industrial economy are therefore generally sheltered from the operation of supply and demand by government regulations of various types, price supports, or tariff protection. Though some attempts have been made to control world commodity markets, these are generally more talk than performance. Some nations, Australia for example, have been able to make enough profit from primary commodity exports to attract capital into the development of industry; but most of the so-called developing countries find their export earnings insecure and insufficient. Their spokesmen complain that the world market system operates in favour of the industrialized nations. (J.Ro.)

FUTURES MARKETS

From very early times, and in many lines of trade, buyers and sellers have found it advantageous to enter into contracts—termed futures contracts—calling for delivery of a commodity at a later date. Dutch whalers in the 16th century entered into forward sales contracts before sailing, partly to finance their voyage and partly to get a better price for their product. From early times, U.S. potato growers in Maine made forward sales of potatoes at planting time. The European futures markets arose out of import trade. Cotton importers in Liverpool, for example, entered forward contracts with U.S. exporters from about 1840. With the introduction of the fast transatlantic Cunard mail services, it became possible for cotton exporters in the United States to send samples to Liverpool in advance of the slow cargo ships, which carried the bulk of the cotton. Futures trading within the United States in the form of “to arrive” contracts appears to have commenced before the railroad days (1850s) in Chicago. Merchants in Chicago who bought wheat from outlying territories were not sure of the arrival time and quality of a delivery. The introduction of “to arrive” contracts enabled the sellers to get a better price for their product and buyers to avoid serious price risk.

Futures trading of this sort in grains, coffee, cotton, and oilseeds also arose in other centres such as Antwerp, Amsterdam, Bremen, Le Havre, Alexandria, and Ōsaka between the 17th and the middle of the 19th centuries. In the process of evolution, “to arrive” contracts became standardized with respect to grade and delivery period, with allowances for grade adjustment when the delivered grade happened to be different. These developments helped to enlarge the volume of trade, encouraging more trading by merchants who dealt in the physical commodity and also the entry of speculators, who were interested not in the commodity itself but in the favourable movement of its price in order to make profits. The larger volume of trading lowered the transaction costs, and by stages the trading became impersonal. The rise of the clearinghouse depersonalized the buyer-seller relations completely, giving rise to the present form of futures trading.

Economic functions of the futures contract. Commodity futures markets provide insurance opportunities to merchants and processors against the risk of price fluctuation. In the case of a trader, an adverse price change brought by either supply or demand change affects the total value of his commitments; and the larger the value of his inventory, the larger the risk to which he is exposed. The futures

Origins
and
develop-
ment

The
kibbutzim
in Israel

The
behaviour
of primary
commodity
markets

market provides a mechanism for the trader to lower the per unit inventory risk on his commitments in the cash market (where actual physical delivery of the commodity must eventually be made) through what is known as hedging. A trader is termed a hedger if his commitments in the cash market are offset by opposite commitments in the futures market. An example would be that of a grain elevator operator who buys wheat in the country and at the same time sells a futures contract for the same quantity of wheat. When his wheat is delivered later to the terminal market or to the processor in a normal market, he buys back his futures contract. Any change of price that occurred during the interval should have been cancelled out by mutually compensatory movements in his cash and futures holdings. The hedger thus hopes to protect himself against loss resulting from price changes by transferring the risk to a speculator who relies upon his skill in forecasting price movements.

For a better understanding of the process involved, the distinctive features of the cash market and the futures market should be made clear. The cash market may be either a spot market concerned with immediate physical delivery of the specified commodity or a forward market, where the delivery of the specified commodity is made at some later date. Futures markets, on the other hand, generally permit trading in a number of grades of the commodity to protect hedger sellers from being "cornered" by speculator buyers who might otherwise insist on delivery of a particular grade whose stocks are small. Since a number of alternative grades can be tendered, the futures market is not suitable for the acquisition of the physical commodity. For this reason, physical delivery of the commodities in fulfillment of the futures contract generally does not take place, and the contract is usually settled between buyers and sellers by paying the difference between the buying and selling price. Several futures contracts in a commodity are traded during a year. Thus, five wheat contracts, July, September, December, March, and May, and six soybean contracts, September, November, January, March, May and July are traded on the Board of Trade of the City of Chicago. The length of these contracts is for a period of about 10 months, and a contract for "September wheat" or "September soybean" indicates the month the contract matures.

Hedging as insurance

Though hedging is a form of insurance, it seldom provides perfect protection. The insurance is based on the fact that the cash and futures prices move together and are well correlated. The price spread between the cash and futures, however, is not invariant. The hedgers, therefore, run the risk that the price spread, known as the "basis," could move against them. The possibility of such an unfavourable movement in the basis is known as basis risk. Thus hedgers, through their commitment in the futures market, substitute basis risk for the price risk they would have taken in carrying unhedged stocks. It must be emphasized, however, that risk reduction is not the final objective with merchants and processors; what they seek to do is to maximize profits.

The availability of capital for financing the holding of inventories depends on whether they are hedged or not. The bankers' willingness to finance them increases with the proportion of the inventory that is hedged. For example, the banks may advance loans to the extent of only 50 percent of the value of unhedged inventories and 90 percent if they are all hedged, a difference explained by the fact that hedging reduces the risk on which the amount of the loan and the interest rate depend. Merchants and processors can therefore derive a twofold advantage from futures trading; they can insure against price decline and they can secure larger and cheaper loans from the banks.

The theory and practice of hedging. There are two rival hypotheses concerning the motives for and costs of hedging. The first of these, advanced by John Maynard Keynes and J.R. Hicks, suggests that risk reduction is the prime motive for hedging and that hedgers pay a risk premium to speculators for assuming risk. The Keynes-Hicks hypothesis states that under normal conditions in commodity markets, when demand, supply, and spot prices are expected to remain unchanged for some months to

come and there is uncertainty in traders' minds regarding these expectations, the futures price, say, for one month's delivery is bound to be below the spot price that traders expect to prevail one month later. This condition exists because inventory holders would be ready to hedge themselves from the risk of price fluctuations by selling futures to speculators below the expected spot price. By selling futures below the expected spot price, according to the theory, inventory holders who hedge pay a risk premium to speculators.

The rival hypothesis of Holbrook Working maintains that hedging is done with the expectation of a profit from a favourable change in the spot-futures price relation, to simplify business decisions, and to cut costs, and not for the sake of reducing risk alone. Hedgers, according to Working, are arbitrageurs; *i.e.*, they take advantage of a temporary price difference between two markets to buy in one and sell in the other. They thus speculate on the basis and assume risk.

A compromise between these rival theories and a more balanced view regarding the need for hedging and the scope of hedging activities is that hedging is motivated by the desire to reduce risks, as suggested by the Keynes-Hicks theory, but that the levels of inventory held by merchants and processors are determined by expected hedging profits, as Working has emphasized.

There are two categories of hedgers in the futures market: they are called short and long hedgers. Short hedgers are merchants and processors who acquire inventories of the commodity in the spot market and who simultaneously sell an equivalent amount or less in the futures market. The hedgers in this case are said to be long on their spot transaction and short in the futures transaction. Wheat merchants or wheat flour mills who either have 100,000 bushels of wheat as inventory or have bought it for later delivery are said to short hedge if they sell 100,000 bushels of wheat in futures contracts. By holding inventories, both merchants and processors can make their purchases when it is most opportune and lower their transaction costs through fewer transactions. Another advantage to the processing firm in holding inventory is that it makes it possible to avoid interruption in production. It must be borne in mind that short hedgers do not normally deliver the physical commodity in fulfillment of the futures contract. They "lift the hedge" by repurchasing the futures contract at the prevailing futures price when they sell the raw material or the processed good in the spot market.

The merchants and processors do not generally hedge all their inventories for the sake of reduced risk. The decision on what part of inventories to hedge is based on their expectations relating to return from holding hedged and unhedged inventories in storage, given the cost incurred in both forms of inventory holding. The return per unit inventory to merchants and processors on their hedged inventories, when liquidated, is the change in the spot price less the change in the futures price and the storage costs. Their return on per unit unhedged inventory is the change in spot price less storage costs.

Long hedgers, in contrast, are merchants and processors who have made formal commitments to deliver a specified quantity of raw material or processed goods at a later date at a price currently agreed upon and who do not now have the stocks of the raw material necessary to fulfill their forward commitment. The parties who have made the commitment generally seek to hedge against the risk of price rise in the raw material between the time of making the forward contract and the time of acquiring the raw material stocks for fulfilling the contract. The hedging is done by buying futures contracts of the raw material equal in quantity to what is needed to fulfill the forward commitment.

The question arises under what circumstances the long hedger might prefer the purchase of futures to the alternative of immediately buying the raw material through spot or forward purchase to meet the obligations of his forward sale. He may prefer buying futures to buying in the cash market (spot or forward) if current cash prices are high because of scarcity. Generally there is an increase in the amount of long hedging when, as the season advances,

The long and short hedger

spot prices rise, inventory holdings fall, and the new crop is not yet available. Long hedging is not as risk-reducing as it may appear at first sight. The long hedger processor, for example, who buys raw material futures to satisfy his forward commitment of the processed good may find that the raw material delivered to him in futures is not of suitable grade and quality to meet the obligations of the forward sale. Quite often, therefore, he may sell his futures contract and purchase raw material of the grade needed. If the spot price of the raw material moves unfavourably relative to the price of the processed good sold forward by him, the long hedger actually increases the risk by buying futures instead of buying the raw material in the cash market. Long hedging, unlike short hedging, may serve to increase risk, and the total risk on long hedging increases with the size of the commitment.

The volume of short hedging tends to be large when stocks in commercial hands are large and when the cash price is below the futures price; a reversal in this situation brings a decline. Conversely, the volume of long hedging is large when stocks are small and the cash price is above the futures price. Short hedging has a marked seasonal pattern, reaching a peak when commercial stocks are largest and the basis is favourable and then declining as the season advances. The seasonal pattern is less marked in long hedging. Generally there is an excess of short over long hedging during the bulk of the crop year.

Apart from hedgers, the futures market includes speculators, and these can also be classified in two categories, namely, long and short speculators. The long speculators are those who expect the price to rise above the current level and assume risks by purchasing futures contracts. Short speculators are those who expect the price to fall. They sell futures contracts. In a futures market the total short selling position, made up of short hedgers and short speculators, and the total long buying position, made up of long hedgers and long speculators, must always be equal. Any excess of short over long hedging must be balanced by an equal excess of long over short speculation. Since short hedging exceeds long hedging for most of the crop year, hedgers are generally short and speculators, therefore, are generally long.

Futures markets have flourished and become important in commodities where sizeable inventories have to be stored and carried forward for meeting the consumption needs of the entire season. Successful futures trading requires a large volume with low transaction costs and that spot and futures prices be well correlated in order to make hedging effective.

Important futures markets. Based on the number and volume of commodities in which active futures trading exists, the United States occupies first place. The Chicago Board of Trade, the largest of the world's futures markets in terms of volume and value of business, is the centre for trading in wheat, corn, oats, rye, soybeans, soybean oil, and soybean meal. About 30 commodities in all are traded on organized exchanges in the United States. The wheat market in Minneapolis, the cotton and wool markets in New York City, and the markets in frozen pork bellies and live hogs in the midwestern United States are among them. The number of commodities in which futures trading takes place are far fewer outside the United States.

There are futures markets for wool in London, Paris, and Sydney; for cotton in Liverpool and Bombay; for sugar in London and Paris; for jute goods in Calcutta; for black pepper in Cochin, India; and for turmeric in Sāngli, India. As a result of government controls on futures markets and also of international commodity agreements, the volume of futures trading in several countries has been adversely affected. The commodity markets in Europe, with few exceptions, have been dormant since the end of World War II. Many of the Indian commodity markets, such as those in gur, jute, and oilseeds, which were once active, have met the same fate. The recurrent arguments in the United States, India, and elsewhere against the futures markets are that they encourage speculation and that the participation of speculators causes price instability. These arguments have led to the demand that markets be controlled or prohibited from functioning. To refute such allegations

requires a comparison between price variations in the presence and absence of speculation, which is impossible for commodities that have futures markets, since it is not meaningful to say for these markets what the price would have been in the absence of speculation. (L.S.V.)

Markets for manufactures

The market for manufactured goods is what economists call "imperfect," because each company has its own style, its own reputation, and its own locations; and all of the arts of advertisement and salesmanship are devoted to making it even more imperfect by attracting buyers to particular brand names. Even small businesses that depend upon outside channels of retail distribution may have the final say in what prices they will charge, and great corporations can differentiate their goods in order to create demand for them.

In this type of market, supply normally is very elastic—that is, responsive to demand—in the short run. Stocks or inventories are held at some point in the chain of distribution; while stocks are running down or building up, there is time to change the level of production, and once a price has been set, it is rarely altered in response to moderate changes in demand. Even in a deep slump, defensive rings may often be formed to prevent price cutting.

In the long run, as well as in the short, supply is responsive to demand in the market for manufactures. It is easier to change the composition of a firm's output than it is to change the production of a mine or a plantation. And when changes in demand are not too rapid, gross profits from one plant can be siphoned off and invested in something quite different. When business is good, moreover, there is continual new investment so that productive capacity is adapted to meeting changing requirements. Workers themselves may not even be aware of changes in the final commodities to which their work contributes, and the level of wages for any grade of factory labour is very little affected by the fortunes of a particular market.

Financial markets

Along with a growth of trade in goods, there has been a proliferation of financial markets. A stock exchange is an organized market for dealing in the securities of businesses and governments. Currently, dealings in securities that came into being in the past predominate over dealings in new issues, and the greater part of industrial investment is financed by retention of profits. Instead of serving mainly as a channel for lending to industry, the chief functions of a stock exchange are to provide a convenient way for owners of property, inherited or newly saved, to place their wealth in income-yielding form, and to provide liquidity and security for them by making a market in which financial assets can readily be realized or redeployed. Dealings in a commodity market, as noted earlier, necessarily contain a speculative element; but in the market for securities, speculation overshadows all else. The art of speculation, as Keynes said, is "to anticipate what average opinion expects average opinion to be." The notorious instability of stock exchanges, with their disturbing effects upon the availability of finance for business and upon the nominal wealth of shareholders, may have repercussions in the market for real goods and services.

Another kind of financial market is the so-called money market. The money market is not an organized exchange; nor is it confined to money in the ordinary sense. The phrase refers to what might be called wholesale transactions in money and short-term credit carried on mainly by large commercial banks. Through the banking system, government authorities (say, the Bank of England or the United States Federal Reserve Board) have some power to control the market for money. But the more the authorities attempt to exercise their powers, the faster the money market develops new organs for bringing lenders and borrowers together and new types of transferrable obligations that are "almost money," in order to escape from control. (An example of this is the huge Eurocurrency and Eurobond markets in which money and financial instruments

The
specula-
tor's role

Chicago
Board of
Trade

Stock
markets
and money
markets

are bought and sold independently of national monetary authorities.) (J.Ro.)

SECURITIES TRADING

Securities are written evidences of ownership giving their holders the right to demand and receive property not in their possession. The most common types of securities are stocks and bonds, of which there are many particular kinds designed to meet specialized needs. This article deals mainly with the buying and selling of securities issued by private corporations. (The securities issued by governments are discussed in the article GOVERNMENT FINANCE.)

Types of corporate securities. Corporations create two kinds of securities: bonds, representing debt, and stocks, representing ownership or equity interest in their operations. (In Great Britain, the term stock ordinarily refers to a loan, whereas the equity segment is called a share.)

Bonds. The bond, as a debt instrument, represents the promise of a corporation to pay a fixed sum at a specified maturity date, and interest at regular intervals until then. Bonds may be registered in the names of designated parties, as payees, though more often, in order to facilitate handling, they are made payable to the "bearer." The bondholder usually receives his interest by redeeming attached coupons.

Since it could be difficult for a corporation to pay all of its bonds at one time, it is common practice to pay them gradually through serial maturity dates or through a sinking fund, under which arrangement a specified portion of earnings is regularly set aside and applied to the retirement of the bonds. In addition, bonds frequently may be "called" at the option of the company, so that the corporation can take advantage of declining interest rates by selling new bonds at more favourable terms and using these funds to eliminate older outstanding issues. In order to guarantee the earnings of investors, however, bonds may be noncallable for a specified period, perhaps for five or 10 years, and their redemption price may be made equal to the face amount plus a "premium" amount that declines as the bond approaches its maturity date.

The principal type of bond is a mortgage bond, which represents a claim on specified real property. This protection ordinarily results in the holders' receiving priority treatment in the event that financial difficulties lead to a reorganization. Another type is a collateral trust bond, in which the security consists of intangible property, usually stocks and bonds owned by the corporation. Railroads and other transportation companies sometimes finance the purchase of rolling stock with equipment obligations, in which the security is the rolling stock itself.

Although in the United States the term debentures ordinarily refers to relatively long-term unsecured obligations, in other countries it is used to describe any type of corporate obligation, and "bond" more often refers to loans issued by public authorities.

Corporations have developed hybrid obligations to meet varying circumstances. One of the most important of these is the convertible bond, which can be exchanged for common shares at specified prices that may gradually rise over time. Such a bond may be used as a financing device to obtain funds at a low interest rate during the initial stages of a project, when income is likely to be low, and encourage conversion of the debt to stock as earnings rise. A convertible bond may also prove appealing during periods of market uncertainty, when investors obtain the price protection afforded by the bond segment without materially sacrificing possible gains provided by the stock feature; if the price of such a bond momentarily falls below its common-stock equivalent, persons who seek to profit by differentials in equivalent securities will buy the undervalued bond and sell the overvalued stock, effecting delivery on the stock by borrowing the required number of shares (selling short) and eventually converting the bonds in order to obtain the shares to return to the lender.

Another of the hybrid types is the income bond, which has a fixed maturity but on which interest is paid only if it is earned. These bonds developed in the United States out of railroad reorganizations, when investors holding defaulted bonds were willing to accept an income obligation

in exchange for their own securities because of its bond form; the issuer for his part was less vulnerable to the danger of another bankruptcy because interest on the new income bonds was contingent on earnings.

Still another hybrid form is the linked bond, in which the value of the principal, and sometimes the amount of interest as well, is linked to some standard of value such as commodity prices, a cost of living index, a foreign currency, or a combination of these. Although the principle of linkage is old, bonds of this sort received their major impetus during the inflationary periods after World Wars I and II. In recent years they have had the most use in countries in which the pressures of inflation have been sufficiently strong to deter investors from buying fixed-income obligations.

Stock. Those who provide the risk capital for a corporate venture are given stock, representing their ownership interest in the enterprise. The holder of stock has certain rights that are defined by the charter and bylaws of the corporation as well as by the laws of the country or state in which it is chartered. Typically these include the right to share in dividends and other distributions, to vote for directors and fundamental corporate changes, and to inspect the books of the corporation, and, less frequently, the "pre-emptive right" to subscribe to any new issue of stock. The stockholder's interest is divided into units of participation, called shares.

A stock certificate ordinarily is given as documentary evidence of share ownership. Originally this was its primary function; but as interest in securities grew and the capital market evolved, the role of the certificate gradually changed until it became, as it is now, an important instrument for the transfer of title. In some European countries the stock certificate is commonly held in bearer form and is negotiable without endorsement. To avoid loss, the certificates are likely to be entrusted to commercial banks or a clearing agency that is able to handle much of the transfer function through offsetting transactions and bookkeeping entries. In the United States, certificates usually are registered in the name of the owner or in a "street name"—the name of the owner's broker or bank; the bank may for legal reasons use the name of another person, known as a "nominee." When a certificate is held in the name of a broker or bank nominee, the institution is able to make delivery more readily and the transfer process is facilitated. Investors, for legal or personal reasons, may prefer to keep the certificates in their own names.

A corporation may endow different kinds or classes of stock with different rights. Preferred stock has priority with respect to dividends and, if the corporation is dissolved, to the division of assets. Dividends on preferred stock usually are paid at a fixed rate and are often cumulated in the event the corporation finds it necessary to omit a distribution. In the latter circumstance the full deficiency must be cleared before payments may be made on the common shares. Participating preferred stock, in addition to stipulated dividends, receives a share of whatever earnings are paid to the common stock. Participation is usually resorted to as an inducement to investors when the corporation is financially weak. Although a preferred issue has no maturity date, it may be given redemption terms much like those of a bond, including a conversion privilege and a sinking fund. Preferred stockholders may or may not be allowed to vote equally with common stockholders on some or all propositions or more characteristically may vote only upon the occurrence of some prescribed condition, such as the default of a specified number of dividend payments.

Common stock, in some countries called ordinary shares, represents a residual interest in the earnings and assets of a corporation. Whereas distributions to bonds or preferred stock are ordinarily fixed, dividends paid on common stock are set at the time of payment by the directors and tend to vary with earnings. The market price of common stock is likely to move in a relatively wide range, depending on investors' expectations of earnings in the future.

Options. An option contract is an agreement enabling the holder to buy a security at a fixed price for a limited period of time. One form of option contract is the stock

Calling
bonds

Income
bonds

Warrants

purchase warrant, which entitles the owner to buy shares of common stock at designated prices and according to a prescribed ratio. Warrants are often used to enhance the salability of a senior security, and sometimes as part of the compensation paid to bankers who market new issues.

Another use of the option contract is the employee stock option. This is used to compensate key executives and other employees; it is normally subject to a variety of restrictions and is generally nontransferable. Stock rights, like warrants, are transferrable privileges permitting stockholders to buy another security or a portion thereof at a specified price for an indicated period of time. The stock right allows stockholders to subscribe to additional shares of stock in proportion to their present holdings. Stock rights usually have a shorter life-span than warrants, and their subscription price is below, rather than above, the market price of the common stock.

The marketing of new issues. The marketing of securities is an essential link in the mechanism that transfers capital funds from savers to users. The transfer may involve intermediaries such as savings banks, insurance companies, or investment trusts. The ultimate user of the funds may be a corporation or any of the various levels of government from municipalities to national states.

The growth of public debt throughout the world has made governments increasingly important participants in the markets for new securities. They have had to develop financing techniques with careful attention to their influence on the markets for nongovernmental securities. The treasuries must carefully study interest rates, yield patterns, terms of financing, and the distribution of holdings among investors.

Local governments are usually subject to various statutory restrictions that must be carefully observed when offering a new issue for sale. Local government bonds are distributed through investment bankers who buy them and reoffer them to the public at higher prices and correspondingly lower yields. Sometimes the terms of the offer are negotiated. In the United States, however, a more prevalent means of selling state and local bonds is through competitive bidding, in which the issuer announces a contemplated offering of bonds for a designated amount, with specified maturity dates, and for certain purposes. Syndicates of investment bankers are formed to bid on the issue, and the award is made to the group providing the most favourable terms. The winning syndicate then reoffers the bonds to the public at prices carefully tailored to be competitive with comparable obligations already on the market and to provide a suitable profit margin.

Marketing corporation securities

The financial manager of a company requiring additional funds has a number of alternative courses of action open to him. He may do all of his financing through commercial banks by means of loans and revolving credit arrangements that, in essence, are formalized lines of credit. Or, he may prefer to raise capital through the sale of securities. If he chooses to do the latter, he may undertake a private sale with an institutional investor such as an insurance company, permitting him to avoid both the complicated procedures of a public distribution and the risks of unsettled market conditions. On the other hand, a private placement of this sort deprives the issuer of the favourable publicity flowing from a successful public offering; it may not afford sufficient resources for very large firms with continuing demands for capital; and it involves rather restrictive legal requirements.

A company that elects to float its securities publicly in the capital market will ordinarily utilize the services of an investment banker. The investment banker may buy the securities from the issuer and seek a profit by selling them at a higher price to the public, thereby assuming the market risks. If the issue is large, the originating investment banker may invite other houses to join with him in purchasing the issue from the company, while to facilitate disposal he may form a selling group to take over the issue from the buying firms for resale to the public. In lieu of buying the securities from the issuer, the investment banker may act as an agent and receive a commission on the amount sold. If the issuer negotiates the selection of an investment banker, the banker will serve as financial

counsel and offer advice on the timing and terms of the new issue. If the selection is by competitive bidding, the relationship is likely to be more impersonal.

An accepted principle of modern finance is that investors are entitled to knowledge about the issuer in order to appraise the quality of the securities offered. A number of countries now require issuers to file registration statements and provide written prospectuses.

The security markets of Europe do not have the aggressive investment-banking machinery developed in the United States. European commercial banks, on the other hand, play a much more important role in financing the needs of industry than do commercial banks in the United States and Great Britain.

In the 1960s, a number of industrial nations faced increasing difficulties in meeting their financing needs through local capital markets. Several issuers began to float securities that were payable in any of 17 different European currencies. This marked what might be called the beginning of an "international securities" market. Efforts were also made to issue bonds on a parallel basis in different countries with each portion denominated in the currency of the country in which it was sold. For various legal and technical reasons, these methods did not attract a wide following.

Another factor that hastened the growth of a European securities market was the balance of payments problem confronting the United States in the 1960s. Certain legislative enactments substantially shut the capital market of the U.S. to foreign issuers; and other restraints were imposed on foreign lending by United States financial institutions and on direct foreign investment by United States corporations. As a result, a number of multinational corporations headquartered in the United States were forced to seek financing in overseas securities markets for the expanding business of their foreign subsidiaries. United States and foreign investment bankers joined in syndicates to float these securities. The process was facilitated by the growth of an international market in Eurodollars, representing claims on dollars deposited in European banks. The bulk of the new bonds offered abroad were denominated in Eurodollars.

Eurodollars market

During the period 1957-65, when this new European market came into being, the volume of foreign bonds issued publicly rose from \$492,600,000 to \$1,489,500,000. The principal and most consistent borrowers were in Canada, Australia, Japan, Norway, Israel, Denmark, and New Zealand. In all of these countries, the major borrower was the national government, except in Canada, where the political subdivisions were the major borrowers. In West Germany, Great Britain, and the United States, the only borrowers in international markets were private units.

The machinery of securities trading. Organized securities markets and stock exchanges are a product of economic development. In the early years of economic growth, most of a country's industrial units are small and their capital requirements relatively modest. The rate of saving is low, and institutions for channelling private savings into investment are generally lacking. As the economy progresses and national income grows, new institutions enter the financial picture to direct the mounting volume of savings into productive outlets. The appearance of growing numbers of individual and institutional investors creates a need for trading markets to speed up transactions and enable stockholders swiftly and easily to convert their holdings to cash.

At this stage of development, corporations usually meet less of their financing needs through direct sales of securities in the new issue market and obtain a larger percentage through reinvesting their own earnings. This plowing back of earnings is not insensitive to the judgment of investors: if the prospects of a company are good, investors bid up the price of its shares in the trading market and show a willingness to forego dividends for the possibility of long-term capital gains achieved through internal growth. Thus, when a company is able to finance its expansion by means of reinvested earnings rather than by new stock issues, the trading segment becomes the more important aspect of the capital market.

The development of stock exchanges

History. Stock exchanges grew out of early trading activities in agricultural and other commodities. Traders in European fairs in the Middle Ages found it convenient to use credit, which required the supporting documents of drafts, notes, and bills of exchange. The French stock exchange may be traced as far back as the 12th century, when trading occurred in commercial bills of exchange. To regulate these incipient markets, Philip the Fair (1268–1314) created the profession of *courratier de change*, the forerunner of the modern French stockbroker, or *agent de change*. At about the same period in Bruges, then a prosperous centre of the Low Countries, merchants took to gathering in front of the house of the Van der Buerse family to engage in trading. From this custom, the name of the family became identified with trading, and eventually “bourse” came to signify a stock exchange. From similar roots in trade and commerce, the institutional beginnings of stock exchanges appeared during the 16th and 17th centuries in other great trading centres throughout the world—Amsterdam, Great Britain, Denmark, Germany.

The growth of trade created a need for banks and insurance companies. Political developments caused governments to seek new sources of funds. This combination of expanding activity and intermittent capital shortages stimulated the early issuers of securities—governments, banks, insurance companies, and some joint-stock enterprises, particularly the great trading companies. From the existing exchanges for commercial bills and notes, it was an easy and logical transition to the establishment of stock exchanges for securities. By the early 1600s, shares of the Dutch East India Company were being traded in Amsterdam; in 1773, London stock dealers who had previously been meeting in coffeehouses moved into their own building; and by the 19th century, trading in securities on a formal basis was common in the industrialized nations.

The evolution of stock exchanges continued. In Great Britain, progress has for the most part been internal and voluntary; the London Stock Exchange has regulated its own activities. The French stock exchanges, in contrast, are directly subject to law, and the operations of the *agents de change* have been affected by national decrees. At one time, there were three markets for securities in Paris: an official market called the Parquet (the “floor”); a semi-official market, the Coullisse (the “wing”); and the Hors Côté (the “outside”), an unregulated market in unlisted securities. In 1929, the Hors Côté was subjected to official regulation and in the following decade its activities were absorbed into the Coullisse, which in turn was combined with the Parquet in a reorganization in 1961. In Belgium the exchanges have had a mixed history. Strict governmental controls were imposed in 1801 and not removed until 1867. Following the economic crisis of 1929–34, the pendulum swung the other way, and the exchanges were once more placed under the control of central authority. In Switzerland, the exchanges have been governed by cantonal (state) law.

Historical events have left their mark upon the development of stock exchanges in some countries. Mining, rather than trade and commerce, was the impelling influence in the establishment of stock exchanges in South Africa and Canada. In Germany, the Berlin Stock Exchange lost its dominant role after World War II, and its position was assumed by exchanges in Frankfurt and Düsseldorf. The Japanese securities markets were revolutionized following World War II, when a new securities law was enacted patterned after the U.S. model. A campaign to distribute stock formerly held by the large *zaibatsu* (family-owned combines) and semigovernment corporations greatly increased public stock ownership, which in turn contributed to the considerable growth of trading on the nine Japanese exchanges. Another post-World War II development was the interest of the governments of developing countries in the use of stock exchanges to facilitate external financing.

Securities markets in the United States began with speculative trading in issues of the new government. In 1791, the country's first stock exchange was established in Philadelphia, then the leading city in domestic and foreign trade. An exchange in New York was set up in 1792, when 24 merchants and brokers decided to charge commissions

while acting as agents for other persons and to give preference to each other in their negotiations. They did much of their trading under a tree at 68 Wall Street. Government securities formed the basis of the early trading. Stocks of banks and insurance companies added to the volume of transactions. The building of roads and canals brought more securities to the market. In 1817 the New York brokers decided to organize formally as the New York Stock and Exchange Board. Thereafter, the stock market grew with the industrialization of the country. In 1863, the New York Stock Exchange adopted its present name. During the Civil War additional exchanges were organized, one of them the forerunner of the present American Stock Exchange, the second largest stock market in the country.

Organization. All stock exchanges perform similar functions with respect to the listing, trading, and clearing of securities. They differ in their administrative machinery for handling these functions.

The London Stock Exchange, the largest in the world in terms of the number and variety of domestic and international securities traded, is an independent institution not subject to governmental regulation. It resembles a private club with its own constitution and operating rules, administered by a council that, except for the government broker who is an *ex officio* nonvoting member, is elected by the membership. Operating responsibility is vested in the secretary and his staff.

In the United States, as in Great Britain, the government does not participate directly in the operations of the exchanges. Since 1933, however, Congress has enacted half a dozen measures that in one way or another affect the securities market. The most important are the Securities Act of 1933, which is primarily concerned with the new-issue market, and the Securities Exchange Act of 1934, which is directed toward the trading market. The latter requires that every stock exchange register with the Securities and Exchange Commission (SEC) as a national securities exchange, unless it is exempted because of the limited volume of its transactions, and that it conform to certain rules in its trading practices. This relationship between the stock exchanges and the SEC is peculiar to the United States; it involves a sort of administrative partnership between the exchanges as private associations—but functioning as quasi-public institutions—and the government. The exchanges generally may adopt policies and issue regulations governing their own operations but are subject to SEC intercession in the event that the Commission believes modifications are required in the public interest.

Most European exchanges are also subject to some form of governmental regulation. The Amsterdam Stock Exchange is a private organization, relatively free to regulate the activities of the market, but the Minister of Finance exercises some supervision under existing legislation. The Zürich exchange is governed by a board of elected members that determines general policy and coordinates the work of the exchange's committees. Of these, the Zürich Cantonal Committee, which directs dealings on the floor of the exchange, is chaired by the head of the Finance Department of the State of Zürich. Although the Frankfurt Stock Exchange is under the direct administration of a Board of Governors elected by its own members, the rules of the exchange must be approved by the authorities of the state of Hesse; and the official specialists, each responsible for the trading in certain securities, are appointed by the Minister of Finance in the state of Hesse. In Brussels the Ministry of Finance is involved in the appointment of members to major committees, while a governmental representative is attached to each exchange to supervise the observance of all rules and laws; in addition, the Banking Commission that is nominated by the government has substantial powers over the admission of securities to public trading. The policy-making Exchange Commission of the Paris bourse is headed by the Governor of the Banque de France, while the regular members are chosen by the Ministry of Finance; the *agents de change* who supervise the trading process are semigovernmental officials. New members of the Italian stock exchanges are appointed by the government from a list of candidates as a result of competitive examinations and therefore have some public status.

New York
Stock
Exchange

Since 1953, membership in the New York Stock Exchange has been limited to 1,366. Only individuals may be members, but they may be partners or stockholders of organizations that do business with the public. Their organizations in such cases are known as member firms. The exchange supervises and regulates member firms in what it considers to be the public interest: a majority of the owners must be engaged primarily in the business of brokers or dealers in securities; exchange approval is required of any shareholder with a 5 percent interest in a member corporation; and all principal officers and directors who are active in the firm must be members or allied members of the exchange. An allied member is subject to the rules of the exchange but does not have the right to engage in transactions on the floor.

To become a member, an individual must acquire, with the approval of the board of governors, a "seat" from a present member or from the estate of a deceased one. Before granting approval, the exchange will investigate such matters as the applicant's past record and financial standing; he will also have to pass an examination demonstrating his knowledge of the securities field.

There are several kinds of brokers on the floor of the exchange. They include the commission broker who executes customer orders placed at or near the current market price; the specialist in one or more issues who, as a broker, executes limited orders for other members and, as a dealer, buys and sells securities for his own account; floor brokers, or "two-dollar" brokers, who execute orders for other brokers at a commission but have no contact with the public; brokers associated with odd-lot firms, who undertake to buy or sell in quantities other than the standard 100-share lot; and "registered traders" who buy and sell for their own account.

To become a member of the London Stock Exchange, an individual must acquire a "nomination" from a retiring member at a price that varies with demand and supply. Every applicant must be approved by at least three-quarters of the stock exchange council. A member may be a broker, dealing as an agent of the public, or a jobber, dealing for his own account with other brokers or jobbers.

Membership in the Paris stock exchange is limited to 85 *agents de change* who supervise activity while the actual work of trading and executing orders is done by their employees and those of the exchange. To become an *agent de change*, an applicant must meet prescribed standards of education and experience as well as pass a written examination. He must be nominated by a retiring member or the heirs of a deceased member and make a deposit guaranty. He is formally appointed by the Minister of Finance.

Other European exchanges set eligibility requirements of character, experience, and financial standing, and some have educational requirements as well. In Brussels, in addition to the completion of six consecutive years in a broker's office, a candidate must have a degree in commercial science or economics and pass a professional examination. In Germany, Switzerland, and Sweden, the brokerage business is dominated by banks.

Members of Japanese exchanges must be corporations doing a business in securities. There are two kinds of members: regular members, who buy and sell for customers or for their own accounts, and *saitori*, who act principally as intermediaries for regular members.

Trading procedures. Most stock exchanges are auction markets, in which prices are determined by competitive bidding. In very large, active markets, the auction is continuous, occurring throughout the day's trading session and for any security in which there is buying and selling interest. In smaller markets the names of the listed stocks may be submitted in some form of rotation, with the auction occurring at that time; this process is described as a "call market."

Trading methods on all the exchanges in the United States are similar. In a typical transaction for a security listed on the New York Stock Exchange, a customer gives an order to an employee in a branch or correspondent office of a member firm, who transmits it either indirectly through the firm's New York office or, as is becoming increasingly common, directly to a receiving clerk on the

floor of the exchange. The receiving clerk summons the firm's floor broker, who takes the order, goes to the post where the stock is traded, and participates in an auction procedure as either buyer or seller. If the order is not a market order calling for immediate action, the broker turns it over to an appropriate specialist who will execute it when an indicated price is reached.

As in any auction market, securities are sold to the broker bidding the highest price and bought from the broker offering the lowest price. Since the market is continuous, buyers and sellers are constantly competing with each other. In the New York Stock Exchange, the specialist plays an important role. As a principal, he has the responsibility of buying and selling for his own account, thereby providing a stabilizing influence; as an agent, he represents other brokers on both sides of the market when they have orders at prices that cannot be readily executed.

With the growing demand for stocks on the part of institutions such as insurance companies, mutual funds, pension funds, and so forth, the size of orders consumed on the New York Stock Exchange has grown. The common way of handling these big blocks on the floor of the exchange has been to break them into smaller orders executed over a period of time. Another method is to assemble matching orders in advance and then "cross" them, executing the purchase or sale at current prices in accordance with prescribed rules; since the broker initially may have obtained the matching orders off the floor, this procedure assumes some of the aspects of a negotiated rather than a pure auction transaction. It is only a step from this to so-called block positioning, in which the broker functions as a principal and actually buys the block from the seller and distributes the securities over a period of time on the floor of the exchange.

When none of these methods appears feasible, the exchange permits certain special procedures. A *secondary distribution* of stock resembles the underwriting of a new issue, the block being handled by a selling group or syndicate off the floor after trading hours, at a price regulated by the exchange. In an *exchange distribution* a member firm accumulates the necessary buy orders and then crosses them on the floor. This is distinguished from an ordinary "cross" because the selling broker may provide extra compensation to his own registered representatives and to other participating firms. A *special offering* is the offering of a block through the facilities of the exchange at a price not in excess of the last sale or the current offer, whichever is lower, but not below the current bid unless special permission is obtained. The terms of the offer are flashed on the tape. The offerer agrees to pay a special commission. A *specialist block purchase* permits the specialist to buy a block outside the regular market procedure, at a price that is somewhat below the current bid.

Trading on the London Stock Exchange is carried on through a unique system of brokers and jobbers. A broker acts as an agent for his customers; a jobber, or dealer, transacts business on the floor of the exchange but does not deal with the public. A customer gives an order to a brokerage house, which relays it to the floor for execution. The receiving broker goes to the area where the security is traded and seeks a jobber stationed in the vicinity who specializes in the particular issue. The jobber serves only in the capacity of a principal, buying and selling for his own account and dealing only with brokers or other jobbers. The broker asks the jobber's current prices without revealing whether he is interested in buying or selling. The broker may seek to narrow the spread between the bid and ask quotations or he may approach another jobber handling the same issue and undertake the same bargaining process. Eventually, when satisfied that he has obtained the best possible price for his client, the broker will complete the bargain.

A broker is compensated by the commission received from the customer. The jobber seeks to maximize his profitable business by adjusting his buying and selling prices. As the ultimate dealer in the London market, the jobber's activities provide a stabilizing factor, but unlike the specialist on the New York Stock Exchange, the jobber is under no obligation to help support prices. The grow-

ing importance of institutional customers has increased the size of transactions in the London market as it has in the U.S., and therefore the jobber has been compelled to risk larger sums. To offset this risk, arrangements for a particularly large order may be negotiated beforehand and the transaction put through the floor as a matter of procedure, with the jobber accepting a minimum "turn." Although the jobbing system provides a continuous market, it does not employ the auction bidding of the New York exchange.

The trading procedures of other major exchanges throughout the world employ the principles that have been described above, although they vary in their application of them. In the exchanges of Paris, Brussels, Copenhagen, Stockholm, and Zürich, some form of auction system is employed: prices are established through bids and offers made on specific securities at particular periods of time. In Tokyo, trading is continuous and orders are consummated through the *saitori* members, who keep order sheets on all transactions. Unlike the specialist on the New York exchange or the jobber in London, however, the *saitori* does not trade for his own account but serves only as an intermediary between regular members. In Amsterdam, trading in active securities is done directly between members during designated trading periods: specialists function as intermediaries between buyers and sellers.

Types of orders. The simplest method of buying stock is through the *market order*. This is an order to buy or sell a stated amount of a security at the most advantageous price obtainable after the order reaches the trading floor. A *limit* (or *limited*) *order* is an order to buy or sell a stated amount of a security when it reaches a specified price or a better one if it is obtainable after the order comes to the trading floor. In the Amsterdam market, the device of the "middle price" is used: an investor who gives a limit order before the opening will have it executed at the day's median level, or at a price that is better than the limit, whichever is found to be more advantageous to the client.

There are other more specialized types of orders. A *stop order* or *stop-loss order* is an order to purchase or sell a security after a designated price is reached or passed, when it then becomes a market order. It differs from the limit order in that it is designed to protect the customer from market reversals; the stop price is not necessarily the price at which the order will be executed, particularly if the market is changing rapidly. This type of order does not lend itself to the London jobbing system.

An important method of trading in stock is through the buying and selling of options. The most common option contracts are puts and calls. A *put* is a contract that permits the holder to deliver to the purchaser a specified number of shares of stock at a fixed price within a designated period of time, say six months; a *call* entitles him to buy shares from the seller within a given period. For example, a person who buys a stock hoping to sell it later at a higher price may also buy a put as a hedge against a fall in price. The put enables him to sell the stock at the price for which he bought it. If the stock rises he need not use the option and loses only the price of its purchase. Option trading is common in Brussels, Paris, London, and the United States.

In the early days of securities trading, stocks and bonds were often bought at private banking houses in the same way that commodities might be purchased over the counter of a general store. This was the origin of the term "over-the-counter." It is used today to mean all securities transactions that are handled outside the exchanges. Increasingly, this market is being subjected to regulation. The extent and nature of the over-the-counter market varies throughout the world. In the United Kingdom, there is no over-the-counter market as such. In the Netherlands, transactions are illegal if they do not involve a member of the Amsterdam exchange or one of its provincial branches as an intermediary, except with the permission of the Ministry of Finance. On the Paris bourse, one post is provided for trading in unlisted issues. In Belgium, the stock exchange committee organizes, at least once a month, public sales of stocks that are not officially quoted. In Japan a second security section has been introduced

into the major exchanges to provide more effective trading procedures for over-the-counter transactions.

In the United States, the over-the-counter market includes most federal, state, and municipal issues as well as a large variety of corporate stocks and bonds. The National Quotation Bureau, which compiles over-the-counter prices, has furnished quotations on approximately 26,000 over-the-counter stocks. In early 1971, a major development occurred with the introduction of current, computerized quotations on a number of active stocks.

Transactions in the over-the-counter market are executed through a large number of broker-dealers with a complex network of private wires and telephone lines. Their operations are subject to the rules of the National Association of Securities Dealers, Inc., a self-regulating body created in 1939. In 1964 the Congress extended to the larger over-the-counter companies the same requirements as to periodic reporting, proxy solicitation, and insider trading that are applied to dealers in listed stocks.

The over-the-counter market is a negotiated market, as distinguished from the auction markets for listed securities. An investor desiring to trade an over-the-counter security gives his order to a broker functioning as a retailer, who ordinarily shops among various firms to obtain the best possible price.

Because of the difficulty that institutions often experience in disposing of large blocks of listed securities on the exchanges, nonmember firms have set up over-the-counter markets in these issues—principally in those listed on the New York Stock Exchange. Although such transactions are conducted within the framework of the over-the-counter market, their prices are tied to those on the Exchange. Accordingly, this form of trading has been labelled the "third market." There is now also a "fourth market," consisting of direct transactions between investors without an intermediary. This market also had its origins in the need of the institutions to find ways of executing large transactions. Impetus to such direct dealings has been given by the development of computerized systems to bring together large traders.

The structure of demand for securities. Interest in the ownership of securities has increased greatly in recent decades. Inflationary tendencies have directed attention to stocks as a means of offsetting rising prices. Stock exchanges have cultivated investors through public relations programs. Government regulation has strengthened public confidence in stock trading procedures. Many governments have given support to the capital markets in order to facilitate business financing.

In the United States, in addition to the millions of individuals who own shares of publicly held corporations, many others own shares indirectly through institutions that are large holders of stock, such as investment companies and pension funds. It is difficult to obtain figures for other countries. A major problem of developing countries has been the absence of investors able and willing to buy shares. Many investors in these countries have preferred to place their funds in tangible assets such as land.

Institutions such as insurance companies, mutual funds, pension funds, foundations, and universities have grown very important in the security markets of the United States. Because of their financial responsibilities to others, these institutions have characteristically followed conservative investment policies stressing the purchase of fixed-income securities. The long-term trend toward inflation, along with mounting stock prices, has led the institutions to look more favourably upon common stocks.

Among the most rapidly growing institutions are the mutual funds. Technically, these are known as open-end investment companies because the number of their shares outstanding constantly changes as new shares are sold to investors and old ones redeemed.

Dynamics of the securities markets. Over a long-term period, the movements of stock prices and of general business indicators tend to parallel each other. In studies of business cycles, it has been found that stock prices tend to reach their cyclical peaks and troughs somewhat ahead of general business indicators, and these are therefore generally classified as "leading indicators."

The distribution of stock ownership

Unorganized, or over-the-counter, markets

Long-term and short-term movements of stock prices

The price of a stock reflects the present value of expected future earnings, and the profits of a firm are strongly influenced by the general level of economic activity. The tendency of stocks to lead business may be attributable to investors' preoccupation with the future. Over the years, the trend of stock prices has been upward; since World War II the upward cycles has tended to be of longer duration, while declines have been relatively shorter. Within the generally expansionary movement, changes in share values of specific companies have been mixed, some showing striking long-term gains while others have suffered losses.

Dramatic events sometimes have a special influence on the psychology of investors, driving stock prices down despite improving business conditions. For example, between the fall of 1940 and the spring of 1942, the period immediately prior to and after the entry of the United States into World War II, U.S. stock prices dropped swiftly despite a continued revival of economic activity. In other cases the reasons for a fall in stock prices are not easy to discover.

Stock prices also experience daily changes of substantial size. Technical market analysts attempt to predict these changes by studying patterns in stock prices. Many theorists, however, claim that in a highly competitive market, prices fluctuate primarily as a result of new information that is not likely to appear in any organized fashion; they maintain that successive price movements are independent of each other and take place in a random fashion.

As investors' expectations change over time, their attitudes toward different types of stock change. Buoyant investors lean toward growth stocks, the value of which is expected to increase rapidly; when uncertainty prevails, the preference is for more conservative issues with stable records of earnings. Within any given period, investors' choices of particular stocks vary with their judgments of the related companies.

(S.M.R.)

MONEY MARKET

The functions of a money market

Every country with a monetary system of its own has to have some kind of market in which dealers in bills, notes, and other forms of short-term credit can buy and sell. The "money market" is a set of institutions or arrangements for handling what might be called wholesale transactions in money and short-term credit. The need for such facilities arises in much the same way that a similar need does in connection with the distribution of any of the products of a diversified economy to their final users at the retail level. If the retailer is to provide reasonably adequate service to his customers, he must have active contacts with others who specialize in making or handling bulk quantities of whatever is his stock-in-trade. The money market is made up of specialized facilities of exactly this kind. It exists for the purpose of improving the ability of the retailers of financial services—commercial banks, savings institutions, investment houses, lending agencies, and even governments—to do their job. It has little if any contact with the individuals or firms who maintain accounts with these various retailers or purchase their securities or borrow from them.

The elemental functions of a money market must be performed in any kind of modern economy, even one that is largely planned or socialist, but the arrangements in socialist countries do not ordinarily take the form of a market. Money markets exist in countries that use market processes rather than planned allocations to distribute most of their primary resources among alternative uses. The general distinguishing feature of a money market is that it relies upon open competition among those who are bulk suppliers of funds at any particular time and among those seeking bulk funds, to work out the best practicable distribution of the existing total volume of such funds.

In their market transactions, those with bulk supplies of funds or demands for them, rely on groups of intermediaries who act as brokers or dealers. The characteristics of these middlemen, the services they perform, and their relationship to other parts of the financial mechanism vary widely from country to country. In many countries there is no single meeting place where the middlemen

get together, yet in most countries the contacts among all participants are sufficiently open and free to assure each supplier or user of funds that he will get or pay a price that fairly reflects all of the influences (including his own) that are currently affecting the whole supply and the whole demand. In nearly all cases, moreover, the unifying force of competition is reflected at any given moment in a common price (that is, rate of interest) for similar transactions. Continuous fluctuations in the money market rates of interest result from changes in the pressure of available supplies of funds upon the market and in the pull of current demands upon the market.

Banks and the money market. *Commercial banks.* Commercial banks are at the centre of most money markets, as both suppliers and users of funds, and in many markets a few large commercial banks serve also as middlemen. These banks have a unique place because it is their role to furnish an important part of the money supply. In some countries they do this by issuing their own notes, which circulate as part of the hand-to-hand currency. More often, however, it is checking accounts at commercial banks that constitute the major part of the country's money supply. In either case, the outstanding supply of bank money is in continual circulation, and any given bank may at any time have more funds coming in than going out, while at another time the outflow may be the larger. It is through the facilities of the money market that these net excesses and shortages are redistributed, so that the banking system as a whole can at all times provide the means of payment required for carrying on each country's business.

In the course of issuing money the commercial banks also actually create it by expanding their deposits, but they are not at liberty to create all that they may wish whenever they wish, for the total is limited by the volume of bank reserves and by the prevailing ratio between these reserves and bank deposits—a ratio that is set by law, regulation, or custom. The volume of reserves is controlled and varied by the central bank (such as the Bank of England, the Bank of France, or the Federal Reserve System in the U.S.), which is usually a governmental institution, is always charged with governmental duties, and almost invariably carries out a major part of its operations in the money market.

Central banks. The reserves of the commercial banks, which are continually being redistributed through the facilities of the money market, are in fact mainly deposit balances that these commercial banks have on the books of the central bank or notes issued by the central bank, which the commercial banks keep in their own vaults. As the central bank acquires additional assets, it pays for them by crediting depositors' accounts or by issuing its own notes; thus the potential volume of commercial bank reserves is enlarged. With more reserves, the commercial banks can make additional loans or investments, paying for them by entering credits to depositors' accounts on their books. And in that way the money supply is increased. It may be reduced by reversing the sequence. The central bank can sell some of its marketable assets in the money market or in markets closely interrelated with the money market; payment will be made by drawing down some of the commercial bank reserve balances on its books; and with smaller reserves remaining, the commercial banks will have to sell or reduce some of their investments or their loans. That, in turn, results in a shrinkage of the outstanding money supply. Central bank operations of this kind are called open-market operations.

The central bank may also increase bank reserves by making loans to the banks or to such intermediaries as bill dealers or dealers in government securities. Reduction of these loans correspondingly reduces bank reserves. Although the mechanics of these lending procedures vary widely among countries, all have one feature in common: the central bank establishes an interest rate for such borrowing—the bank rate or discount rate—pivotal in the structure of money market rates.

Money market assets may range from those with the highest form of liquidity—deposits at the central bank—through bank deposits to various forms of short-term paper such as treasury bills, dealers' loans, bankers' accep-

Volume of reserves

tances, and commercial paper, and including government securities of longer maturity and other kinds of credit instruments eligible for advances or rediscount at the central bank. Although details vary among countries, the touchstone of any money market asset other than money itself is its closeness—*i.e.*, the degree of its substitutability for money. So long as the institutions making use of a money market regard a particular type of credit instrument as a reasonably close substitute—that is, treat it as “liquid”—and so long as the central bank acquiesces in or approves of this approach, the instrument is in practice a money market asset. Thus no single definition or list can apply to the money markets of all countries nor will the list remain the same through the years in the money market of any given country.

The international money market. Each central bank usually holds some form of reserve that is acceptable in settling international transactions. International monetary reserves are mainly gold, or “money market assets” in some country whose currency is widely used, such as the United States dollar. The monetary laws of all countries provide for the establishment of some kind of parity between their currencies and those of other countries. This parity may be defined either in terms of gold or in relation to a key currency such as the British pound sterling or the United States dollar, which in turn has a fixed parity with gold. A country maintains the “convertibility” of its currency by standing ready to buy and sell gold or other currencies in exchange for its own at prices within a fixed and rather narrow “spread” above or below the “exchange rate” for its own currency that is implied by the declared parity.

Because world trade continually gives rise to various needs for payment in various currencies, an international money market must exist so that traders with an excess of one currency can use it to buy another currency for which they have a need. Within the scope of convertibility arrangements, this trading in currencies is carried out by skilled intermediaries, usually banks or specialized foreign exchange brokers and dealers. Trading in currencies is extensive both for immediate use (“spot”) and for future (“forward”) delivery. Quotations vary according to changes in supply and demand, over the range between the upper and lower buying and selling prices set by official parity. If no parity has been set quotations may fluctuate widely. If a currency is subject to exchange controls, there may be two or more quotations for different uses of the same nominal currency.

Changes in a country’s balance of payments may affect the usefulness or prestige of its currency. A sustained and substantial balance of payments deficit (outpayments larger than inpayments), for example, will result in continuous large increases in the world supply of its currency, possibly leading to some decline in its acceptability abroad and to a loss of international monetary reserves. At the same time, an outward drain may reduce the reserves of the commercial banks (the base for the domestic money supply), unless the central bank takes offsetting action.

Since 1944 most of the countries that have domestic money markets or that play a role in the international money market have been joined together in the International Monetary Fund, which represents a pooling of part of the foreign exchange reserves (including gold) of more than 100 member countries. Drawings on the pool may be made by member countries to meet some of the reserve drains arising from balance of payments deficits and in amounts related to the quota that each has subscribed.

The internal money markets of a surprisingly high proportion of the countries of the world are quite rudimentary. The work of the money market in these countries is done largely by transfers of deposit balances, government securities, or foreign exchange among a few banks and between them and the central bank. But in nearly all such cases there is genuine discontent with the rigidity of these limited facilities and a desire to develop a structure, as well as instruments and procedures, which would provide the open-market attributes of the arrangements that have evolved in the leading countries. Several of the more fully developed money markets are described below.

The U.S. money market. The domestic money market in the United States carries out the largest volume of transactions of any such market in the world; its participants include the most heterogeneous group of financial and nonfinancial concerns to be found in any money market; it permits trading in an unusually wide variety of money substitutes; and it is less centralized geographically than the money market of any other country. Although there has always been a clustering of money market activities in New York City and much of the country’s participation in the international money market centres there, a process of continuous change during the 20th century has produced a genuinely national money market.

By 1935 the financial crises of the Great Depression had resulted in a basic revision of the banking laws. All gold had been withdrawn from internal circulation in 1933 and was henceforth held by the U.S. Treasury for use only in settling net flows of international payments among governments or central banks; its price was raised to \$35 per ounce, and the U.S. dollar became the key currency in an international gold bullion standard. Domestically, the changes included legislative recognition of the primary importance of unified open-market operations by the Federal Reserve System and delegation to the board of governors of the Federal Reserve System of authority to raise or lower the ratios required between reserves and commercial bank deposits. Although about half of the 30,000 separate banks existing in the early 1920s had disappeared by the mid-1930s, the essential character of commercial banking in the U.S. remained that of a “unit” (or single-outlet) banking system in contrast to those of most other countries, which had a small number of large branch-banking organizations.

The unit banking system. This system has led inevitably to striking differences between money market arrangements in the United States and those of other countries. At times, some smaller banks almost inevitably find that the wholesale facilities of the money market cannot provide promptly the funds needed to meet unexpected reserve drains, as deposits move about the country from one bank to another. To provide temporary relief, pending a return flow of funds or more gradual disposal of other liquid assets in the money market, such banks have the privilege, if they are members of the Federal Reserve System, of borrowing for reasonable periods at their own Federal Reserve bank. At times some large banks, which serve as depositories for part of the liquid balances of many of the smaller ones (including those that are not members of the Federal Reserve System) also find that demands converging on them are much greater than expected. These large banks, too, can borrow temporarily at a Federal Reserve bank if other money market facilities are not adequate to their needs. Because these borrowing needs are unavoidably frequent in a vast unit banking system and, as a rule, do not indicate poor management, the discount rate charged by the Federal Reserve banks on such borrowing is not ordinarily put at punitive or severe penalty levels—thus, contrary to practice in many other countries, the central bank does not always maintain its interest rate well above those prevailing on marketable money market instruments. To avoid abuse, there is continuous surveillance of the borrowing banks by the Federal Reserve banks.

Along with this practice of borrowing at a Federal Reserve bank has developed the market for “federal funds.” This specialized part of the money market provides for the direct transfer to a member bank of balances on the books of a Federal Reserve bank in return for payment of a variable rate of interest called the “federal funds rate.” These funds are immediately available. There are transactions, too, in funds that are on deposit at commercial banks—by means of loans between banks, or through loans by one large depositor to another. Because these must be collected through a clearing process, they are usually called “clearinghouse funds.”

Money market instruments. Transactions in federal funds and clearinghouse funds are further supplemented by transactions in which either kind of money is exchanged for some other liquid, money market instrument, most frequently government securities. The magnitude of

the market for government securities became so great after World War II that it overshadowed all other elements of the money market. Trading in outstanding "governments" is virtually all done through dealers who buy and sell for their own account at prices which they quote on request (standing ready to "bid" for or to "offer" any outstanding issue). Most of these dealers have head offices located in New York City, but all are engaged in nationwide operations. Their transactions and the lending arrangements through which they finance their own inventories of government securities have evolved into a particularly sensitive indicator of the pressures of supply and demand on the money market from day to day. The most common form of dealer financing is the repurchase agreement, through which dealers sell parts of their inventory temporarily, subject to repurchase.

Closely interrelated, often through trading operations conducted by the same dealers, are the much smaller markets for bank drafts, bills of exchange, and commercial paper. Alongside these other markets and actually somewhat larger in outstanding volume are the markets for securities issued by various "agencies" created by federal statute, such as the Federal Home Loan banks and Federal Land banks. Another money market instrument is the negotiable time certificate of deposit (CD), issued in large volume by commercial banks, which first became significant in 1962. While the owner of a time CD cannot withdraw his deposit before the maturity date initially agreed upon, he can sell it at any time in a secondary market that is conducted by government securities dealers.

The Federal Reserve System conducts day-to-day operations in the money market on its own initiative in order to assist the smooth working of the nation's financial machinery and to exert a general influence aimed at fostering economic growth and limiting economic instability. Its transactions include substantial outright purchases or sales of government securities, relatively small purchases and run-offs of bankers' acceptances, and a considerable volume of loans made for a few days at a time to dealers in government securities or acceptances in the form of repurchase agreements. While it is still the commercial banks as a group that have the greatest continuing need for the combined facilities of the nationwide money market, there is frequent and continuous participation by a great variety of institutional investors who channel the public's savings into various uses and who must always also make some provision for their own liquidity.

Perhaps the most unusual feature in the composition of the U.S. money market is the great importance attained by nonfinancial business concerns and local units of government since World War II. Corporate treasurers and the treasurers of many states and local political subdivisions and authorities have become so keenly sensitive to the profitable possibilities of managing their own liquid holdings instead of relying on the commercial banks as most had done formerly that this group at times provides nearly as large a part of the volatile financing needs of government securities dealers, for example, as comes from the banks. Moreover, banks outside New York City sometimes supply more of the financing needed by these dealers than do the traditional "money market banks" in New York City. The nationwide character of the money market is also shown by the participation of nearly 200 banks in the federal-funds market—banks that are widely scattered among all Federal Reserve districts, although the bulk of all transactions is executed through facilities located in New York.

While the U.S. money market has become truly national, it still needs a final clearing centre upon which the net impact of changes in overall supply or demand can ultimately converge and where the final balancing adjustments of the market as a whole can be accomplished. In filling that need, New York City continues to be the centre of the national money market.

The British money market. *The discount houses.* In Great Britain the money market consists of a number of linked markets, all of them concentrated in London. The 12 specialist banks known as discount houses have the longest history as money market institutions; they

have their origin in the London bill broker who in the early 19th century made the market in inland commercial bills. By selling bills through this market, the growing industrial and urban areas were able to draw upon the surplus savings of the agricultural areas. Quite early many bill brokers began to borrow money from banks in order to buy and hold these bills, instead of simply acting as brokers, and thus became the first discount houses. Since then the major assets held by the discount houses have at different times been commercial bills (first inland bills as described above and later bills financing international trade), treasury bills, and short-dated government bonds. During the 1960s there was a resurgence of the commercial bill, which finally became the discount houses' largest single class of asset, only to be overtaken later by the certificate of deposit.

Important changes were introduced into the British monetary system in 1971, but money at call with the discount houses retained its role as a reserve asset. Such is the safety and liquidity of call money that, despite the fractionally lower rate on it compared with other reserve assets, the banks hold about half of their required reserves in this form. This in turn provides the discount houses with a large pool of funds, which they invest in relatively short-dated assets, of which the most important is sterling certificates of deposit, followed by commercial bills, local authority securities, and treasury bills. This pattern of assets is greatly influenced by the fact that all call loans to the discount houses are secured loans, parcels of assets being deposited *pro rata* with the lending banks as security, and the assets held by the discount houses must therefore be suitable for use as such security.

They also need to hold a substantial proportion of assets that are rediscountable at the Bank of England in case of need, and the Bank of England limits their holding of assets other than public sector debt to a maximum of 20 times their capital resources.

On the liabilities side of the discount house's balance sheet, operating in call money is part of its day-to-day work. A bank that expects to make net payments to other banks during the day (for example, in settlement of checks paid by its customers to their customers) will probably call in some of its call loans, and by convention this is done before noon. Since the banks that have called in money then pay it to other banks, these other banks will have an equal amount to re-lend to the discount houses in the afternoon. Thus the discount houses can "balance their books"—borrow enough money in the afternoon to replace the loans called from them in the morning.

It is not uncommon for perhaps £100,000,000 to be called from, and re-lent to, the discount houses on an active day.

There is one main reason why the money position may not balance in this way. The British government accounts are kept with the Bank of England, which does not lend at call as other banks do. Thus net payments into or out of these government accounts will cause a net shortage or surplus of money for the discount houses in the afternoon and will tend to cause money rates to rise or fall. The Bank of England can allow such shortages or surpluses to affect interest rates, or it can offset them by buying or selling bills or by lending overnight to the discount houses at market rates. Even if the Bank of England does not act in this way to meet a shortage of funds, the discount houses are always finally able to secure the funds they need by their right to borrow from the Bank of England (the lender of last resort) against approved security at the "minimum lending rate" (the penalty rate).

On the assets side of their balance sheets, the discount houses are active dealers in a number of the assets they hold. They make the market in sterling certificates of deposit and in commercial bills, quoting buying and selling rates for different maturities. They also quote selling rates for treasury bills that they acquire at the weekly tender in competition with each other and with any other banks that may tender, including the Bank of England. Most of these other banks tender for treasury bills in order to hold them to maturity, but the discount houses sell theirs on the average when only a few weeks of the bills' 91-day life has passed. A large proportion of these bills is sold to the

Certificates of deposit

Operations in money

Dealing in assets

clearing banks, which do not tender on their own account.

The Bank of England minimum lending rate is normally determined for each week 0.5–0.75 percent above the average treasury bill rate at the previous Friday's tender. The bank, however, has the power to fix it at a different level if it so wishes, and this has been done.

Other markets. Important changes have also occurred outside the discount market described above; after the mid-1950s there was steady growth in public borrowing by local authorities. This led to an active local authority loan market conducted through a number of brokers, where money can be lent on deposit for a range of maturities from two days up to a year (and indeed for longer periods). Much more rapid was the growth after the mid-1960s of the interbank market, in which banks borrow and lend unsecured for a range of maturities from overnight upward. This market also is conducted through brokers, often firms that also operate in the local authority and other markets; a number of these firms of brokers are subsidiaries of discount houses.

In addition to the markets mentioned, there is the gilt-edged (government bond) market on the stock exchange; short-dated bonds are held by the discount houses and by banks and other money market participants, as are short-dated local authority stocks and local authority "yearling" (very short-dated) bonds. With flexibility of bank deposit rates (at least for deposits of large denomination), both banks and nonbank transactors are faced with a wide and competitive range of sterling money market facilities in London.

Finally, mention should be made of the Eurodollar market, because London is its centre; this is an entrepôt market with a very large volume of business in U.S. dollar balances, conducted through brokers (often the same firms that operate in the sterling markets), and U.K. banks are active participants. However, owing to exchange control there has been little significant interaction between the Eurodollar market and the U.K. domestic money market.

The money markets of other countries. *The Canadian money market.* The Canadian money market was substantially broadened in 1954 with the introduction of day-to-day bank loans against Government of Canada treasury bills and other short-term government and government-guaranteed securities. Treasury bills of 91 days' and 182 days' maturity are issued weekly with the occasional offering of a longer maturity of up to one year. Government of Canada bonds and Government of Canada guaranteed bonds are issued at less regular intervals.

Groups involved in the money market are the following: the government, as the issuer of the securities; the Bank of Canada, acting as issuing agent for the government and as a large holder of market material; the chartered banks, as large holders and as distributors and potential buyers and sellers of bills and bonds at all times; the security dealers, as carriers of inventories and traders in such securities; and the public (mainly provincial and municipal governments and larger corporations), as short-term investors.

Treasury bills are sold by competitive tender in which the Bank of Canada, the chartered banks, and a small number of investment dealers participate. Bonds are normally issued at a price at which the yield is in line with outstanding comparable issues.

The central bank, through its tender at the weekly treasury-bill sale, active manipulation of its own bill and bond portfolio, and regulation of the money supply, has workable instruments for active monetary control. For both banks and qualified dealers, the Bank of Canada acts as lender of last resort. The rate is set slightly above the average rate of the last treasury bill auction to discourage regular borrowing.

The German money market. In what was formerly West Germany, where the money market developed strongly after World War II, transactions have been to a large extent confined to interbank loans. In addition, insurance companies and other nonbank investors are also important lenders of short-term funds. Treasury bills and other short-term bills and notes from government agencies (railways and post) were gaining in importance by the 1960s, whereas in 1955 certain nonmarketable securities (the so-

called equalization claims, created during the 1948 currency reform) held by the Bundesbank were transformed into short-term marketable securities in order to obtain suitable market material for the open-market operations of the Bundesbank. Banks are not used to dealing in short-term government securities between each other. They generally either hold these securities to maturity or resell them to the central bank at its buying rates, so that a true money market has not developed.

The market for commercial paper is of some significance, and dealing in it takes place from time to time between banks, especially in times of tight market conditions. Comprehensive regulations have been given through the Bundesbank about the rediscountability of the several kinds of commercial paper.

The influence of the Bundesbank on the monetary situation through open-market operations by the 1960s was greatly hampered by the vast liquidity of the banking system as a consequence of the persistence of Germany's favourable balance of payments situation.

The French money market. The French money market is fairly well established, but its size is restricted by the fact that in France currency still plays an important role in the money supply, whereas by regulations the nonfinancial private sector of the economy is excluded from dealing in the market. Banks as well as a few public or semipublic agencies working in the financial sphere and intermediaries—brokers and discount houses—constitute the market. Transactions take place in commercial paper and in treasury bills. The monetary authorities maintain a special bookkeeping system for all the treasury bills held by banks and other financial institutions, under which such bills are not represented by actual certificates but by entries in special accounts administered by the Banque de France for the treasury.

The central bank's open-market operations, which were normally limited to smoothing out disturbances in the local money market, have gained importance in recent years. Open-market transactions are effected to keep domestic money market rates in line with international rates, in an effort to prevent unwanted capital flows. The possibilities of the central bank's influencing the monetary situation through the money market are limited to the large government needs for short-term funds, no market for long-term government borrowing being established.

The Japanese money market. In Japan's rapidly growing economy the demand for funds, both short-term and long-term, has been persistently strong. Commercial banks and other financial institutions have therefore had an important role. The monetary authorities (the Ministry of Finance and the Bank of Japan) have been unwilling to allow market forces to equilibrate demand and supply in many financial markets for fear that interest rates would become excessively high. Most interest rates have been set administratively at levels high by international comparison (until the late 1960s) but lower than market forces would have dictated. Monetary policy is implemented by controls on both the availability of credit and its cost.

Under these circumstances, Japan has had a very restricted money market. The market for short-term government securities is negligible; the low, pegged interest rate means that the Bank of Japan is the main buyer and that open-market operations are impossible. Transactions in commercial paper are minimal, being discouraged because they would tend to undermine the structure of interest rates and financial institutions.

Only the call money market is well developed. It is restricted to transactions among financial institutions. The interest rate on call money has been relatively free, and persistently above most other short-term and long-term rates. Although small amounts are lent overnight, most are "unconditional loans" (repayment after one day's notice, with a minimum of two days) or "over-month-end-loans" (repayment on a fixed day the following month). The pattern of flows is rather stable, despite seasonal and cyclical fluctuations. City banks are the major borrowers; they have a strong demand for loans by large enterprises and use call funds as a major source of liquidity. Major lenders are local banks, trust banks, credit associations,

Government of Canada bonds

Call money market in Japan

and agricultural cooperatives, which collect individual urban and rural savings and are attracted by the high yields, liquidity, and low risk of call loans relative to other uses. Call brokers help make a market, though most funds flow directly from one financial institution to another. About three-quarters of the funds flow through the Tokyo market, and there are also call markets in Osaka and Nagoya.

Money markets in developing countries. Well-developed money markets exist in only a few high-income countries. In other countries money markets are narrow, poorly integrated, and in many cases virtually nonexistent. Despite the many differences among countries, one can say in general that the degree of development of a country's financial system, including its money markets, is directly related to the level of its economy. Most very-low-income countries have limited financial systems in which money markets play no role. In many former colonies, notably in Africa, expatriate commercial banks had substituted for a local money market; the banks met fluctuations in loan demand by changing their balances at head offices in London or elsewhere. More recently, government policies have encouraged these banks to develop domestic channels for temporary surpluses and deficits. Persistent inflation has been another factor inhibiting the growth of money markets in developing countries, notably in Latin America.

Most developing countries, except those having socialist systems, have the encouragement of money markets as a policy objective, if only to provide outlets for short-term government securities. At the same time many of these governments pursue low-interest-rate policies in order to reduce the cost of government debt and to encourage investment. Such policies discourage saving and make money market instruments unattractive. Nevertheless, a demand for short-term funds and a supply of them exist in all market-oriented economies. In many developing countries these pressures have led to "unorganized money markets," which are often highly developed in urban areas. Such markets are unorganized because they are outside "normal" financial institutions; they manage to escape government controls over interest rates; but at the same time they do not function very effectively because interest rates are high and contacts between localities and among borrowers and lenders are limited. In all developing countries traditional forms of moneylending continue, particularly for agriculture and small enterprise.

(H.T.P./R.V.R./R.F.G.A./C.G.)

The markets and social welfare

The branch of academic teaching called the "economics of welfare" is nowadays generally expounded by starting with a presumption in favour of *laissez-faire* and free competition and then following it by a list of reservations and exceptions that destroy its validity. One is first asked to imagine a particular group of individuals, each with his tastes for a predetermined list of specific commodities, with a predetermined endowment of labour power, with specified equipment and stocks, and with a given body of technical knowledge. These individuals can produce various alternative combinations of commodities. If resources were fully utilized, it would be impossible to produce more of any one commodity without producing less of others. One could then speak of an "opportunity cost" of each commodity—that is, the cost of producing a little more of it in terms of the amounts of other commodities that would have to be sacrificed.

On the other side of the market, the same individuals appear as consumers. If their tastes are all alike, it can be argued that there is one set of outputs and prices at which the relative subjective valuations are proportional to relative opportunity costs. This situation is the optimum in the sense that any move away from it would reduce the total satisfaction. The argument is not so simple when different consumers have different tastes. At any one pattern of prices established in the market (with full utilization of resources), each consumer can adjust his purchases so as to maximize his satisfaction at those prices. Then, starting from that set of prices, it can be shown that no one consumer could be made better off, by changing the com-

position of output, without making some other worse off. This is a powerful argument for the status quo, wherever it may happen to be. At every other point, there would be a different status quo, with a pattern of prices suiting some consumers better and some worse than the one that was chosen at the beginning of the argument.

The first general objection that is now admitted to this scheme (over and above some technical points within its own terms) is that it is purely static. Actual life is lived in the stream of time. In any economy, tastes, commodities, resources, and technology are continuously changing and modifying each other.

Second, there is a deep-seated inconsistency in the assumptions: when every individual pursues his own advantage, atomistic competition cannot persist, for any group of sellers or of buyers can gain by acting together and sharing the benefits among individual members of the group. Perfect competition, like free trade among nations, can persist only when there is some rule of behaviour that overrides pure self-interest.

Third, it is necessary to consider the distribution of consumption between individuals who make up the market. The purchasing power of each is limited by the receipts that he draws from his sales. His receipts depend on the amount of his original endowment and the price of the product he offers to sell. It is obvious that the distribution of the benefits of the market between individuals is quite arbitrary, and the market optimum is therefore not the same as their welfare optimum. When this objection to the *laissez-faire* principle is admitted, it is often suggested that inequalities might be corrected by a system of bounties and taxes; but this has never been taken seriously as a practical scheme.

A fourth set of objections comes under the heading of differences between private and social costs. The classic example is the smoke nuisance, which imposes costs upon the public that the factory concerned cannot be charged for. Examples of this phenomenon on a devastating scale are now daily coming to notice. At the same time there are many socially beneficial activities for which it is impossible to collect adequate payment from individuals in a society of unequal wealth and income—not only performances of grand opera but the whole of education and the health service.

Thus the doctrine that the free play of individual interests in a competitive market maximizes welfare for society as a whole has been demolished by its own exponents. Yet these notions still have an important influence on the formation of ideology. The rationalization grew out of the 19th-century utilitarian philosophers' hedonistic calculus, which regarded all human life as governed by the pursuit of pleasure and avoidance of pain. This train of thought identifies pleasure with consumption, leaving out of account the whole sphere of conditions in which work is carried on, and finally it sees consumption in terms of spending money. The whole economic basis of life is thus reduced to terms of commercial transactions. As the American critic Thorstein Veblen pointed out,

so great and pervading a force has this habit of pecuniary accountancy become that it extends, often as a matter of course, to many facts which properly have no pecuniary bearing and no pecuniary magnitude, as, for example, works of art, science, scholarship, and religion.

THE POLITICS OF THE MARKET

The doctrines of *laissez-faire* are attractive in many ways. If the economy is a self-regulating mechanism and if economics is a system of scientific laws, then moral and political problems are excluded. Questions of social justice do not arise. The function of government is to be strictly neutral between interested parties. But when people come to recognize that the market, by its very nature, is necessarily a scene of conflicting interests, every element in it becomes a moral and political problem. This is distressing because one can no longer rely upon "principles of economics" to provide safe and simple rules for finding the correct solutions.

The intrusion of politics into economic affairs increased dramatically in all the Western industrial nations after

Static
balance of
interests

Laissez-faire as an ideal

1945. Capitalism took on a new shape, with certain characteristics that distinguish it from all former systems. The reasons for this mutation in the free market economy are numerous. Perhaps the most important was the experience of the great slump of the 1930s. After the war, public opinion, business interests, politicians, and administrators were united in resolving that massive unemployment must not be allowed to recur. Indeed, those who earlier had been the strongest adherents of laissez-faire became the strongest supporters of policies to maintain full employment, arguing that this is the best way to preserve the free enterprise system. Some observers of modern capitalism have maintained that the real purpose of the full-employment policy is to maintain profitability for the great corporations, but it cannot be denied that the results are important for all classes.

The traditional doctrines of the economists are still influential. Even John Maynard Keynes, whose critique of the traditional economics had so much to do with creating the new, gave his blessing to the old conception of the free market. He argued that, if our central controls succeed in establishing full employment, there is no objection to be raised against the orthodox theory of the manner in which private self-interest will determine what in particular is produced, or how the value of the final product will be distributed.

Evidently this view owes more to sentiment than to logic. It is impossible to have a policy designed to maintain effective demand in the abstract. Every policy must have some concrete content. The instruments in the hands of government—monetary policy, the exchange rate, taxation, government expenditure—impinge in specific ways upon specific interests. Taxation may be designed to foster either investment or consumption; if the latter, it may be done either by adhering to the principle of “to him that hath shall be given” (that is, by a proportional reduction in direct tax rates) or by favouring lower-income groups. Government investment cannot be neutral either; there is no criterion for allocating funds in a neutral manner between, say, armaments and hospitals.

While economic affairs grow ever more overtly political, at the same time commercial influences spread into new spheres. Old notions of loyalty, service, and proper behaviour are undermined by commercial principles. Continuous inflation and rapid technical change make it necessary for every group to defend itself against the erosion of its relative income and status. It would appear that the new capitalism, for all its benefits, may have been developing internal contradictions of its own that threaten its future stability.

Meanwhile, many of the new and developing nations find that the opening up of their resources and their markets turns out to be more advantageous for established businesses based on successful capitalist industry than for struggling newcomers in their own countries. The economic performance of the so-called developing countries does not give much support to the doctrine that the free play of market forces can be relied upon to maximize human welfare.

(J.Ro./Ed.)

BIBLIOGRAPHY

General works. GLENN G. MUNN, F.L. GARCIA, and CHARLES J. WOELFEL, *Encyclopedia of Banking and Finance*, 9th ed., rev. and expanded (also published as *The St. James Encyclopedia of Banking & Finance*, 1991), provides comprehensive definitions, many with bibliographies. EDWARD I. ALTMAN and MARY JANE MCKINNEY (eds.), *Handbook of Financial Markets and Institutions*, 6th ed. (1987), is a thorough compilation. Detailed information on a variety of markets is provided in FRANCIS A. LEES and MAXIMO ENG, *International Financial Markets: Development of the Present System and Future Prospects* (1975), a descriptive treatment; CHARLES R. GEISST, *A Guide to the Financial Markets*, 2nd ed. (1989), for the general reader; FRANK J. FABOZZI and FRANK G. ZARB, *Handbook of Financial Markets: Securities, Options, and Futures*, 2nd ed. (1986); and PERRY J. KAUFMAN, *Handbook of Futures Markets: Commodity, Financial, Stock Index, and Options* (1984), including the history, regulation, and mechanics of futures trading. Further discussion of financial futures is found in MARK J. POWERS and MARK G. CASTELINO, *Inside the Financial Futures Markets*, 3rd ed. (1991), an explanation of the exchanges and their func-

tions; and NANCY H. ROTHSTEIN and JAMES M. LITTLE (eds.), *The Handbook of Financial Futures: A Guide for Investors and Professional Financial Managers* (1984), a discussion of the market's development, organization, and regulation.

The market in economic doctrine and history. The first chapter of ADAM SMITH, *An Inquiry into the Nature and Causes of the Wealth of Nations* (1776, reprinted frequently), contains his famous discussion of the division of labour. ALFRED MARSHALL, *Principles of Economics*, 9th ed., 2 vol. (1961), conveys his approach to the market. The development of the general equilibrium approach to markets by Leon Walrus and others is well recounted by JOSEPH A. SCHUMPETER, *A History of Economic Analysis*, ed. by ELIZABETH BOODY SCHUMPETER (1954). The best short introduction to the Keynesian Revolution is by MICHAL KALECKI, *Studies in the Theory of Business Cycles, 1933-1939* (1966; originally published in Polish, 1962). These essays were written before the publication of the great work of JOHN MAYNARD KEYNES, *The General Theory of Employment, Interest, and Money* (1935, reissued 1991). A critical account of the theory of imperfect competition is presented in the preface to JOAN ROBINSON, *The Economics of Imperfect Competition*, 2nd ed. (1969, reissued 1976). A slightly different approach is that of EDWARD HASTINGS CHAMBERLIN, *The Theory of Monopolistic Competition*, 8th ed. (1962). Economics without markets are described in KARL POLANYI, *Primitive, Archaic, and Modern Economies*, ed. by GEORGE DALTON (1968), a collection of essays of great interest and originality. ANDREW SHONFIELD, *Modern Capitalism* (1965, reissued 1978), studies the ways in which various countries have adapted their economic administration to modern requirements. A more critical view of modern capitalism is that of JOHN KENNETH GALBRAITH, *The New Industrial State*, 4th ed. (1985). A Marxist view is set forth by PAUL BARAN and PAUL SWEETZ, *Monopoly Capital* (1966). A summary of the attempts at economic reform in the then-existent Soviet Union and other countries with socialist economies is given in MICHAEL ELLMAN, *Economic Reform in the Soviet Union* (1969). The economic problems of the poor countries are examined in GUNNAR MYRDAL, *The Challenge of World Poverty* (1970), a continuation of his monumental work *Asian Drama: An Inquiry into the Poverty of Nations*, 3 vol. (1968), also available in an abridged edition (1971).

Commodity and futures markets. J.R. HICKS, *Value and Capital*, 2nd ed. (1946, reissued 1975), contains a short discussion of the theory of “normal backwardation,” also known as the Keynes-Hicks hypothesis. HOLBROOK WORKING, “New Concepts Concerning Futures, Markets, and Prices,” *American Economic Review*, 52(1):431-459 (June 1962), presents his theories on the role of hedging and the functions of futures markets. The instability of markets for primary commodities is the theme of the Argentinian economist RAUL PREBISCH, *Towards a New Trade Policy for Development* (1964). Recent work on the specific methods and rules of the major world exchanges is included in JOHN BUCKLEY (ed.), *Guide to World Commodity Markets*, 5th ed. (1986). Basic works on commodity futures trading include GERALD GOLD, *Modern Commodity Futures Trading*, 7th rev. ed. (1975); and BRUCE G. GOULD, *The Dow Jones-Irwin Guide to Commodities Trading*, rev. ed. (1981). See also BARRY P. BOSWORTH and ROBERT Z. LAWRENCE, *Commodity Prices and the New Inflation* (1982), comparing the experiences of the United States, West Germany, and Japan, and offering policy proposals; and ROBERT L. ROTHSTEIN, *Global Bargaining* (1979), utilizing the 1974-77 United Nations conference on trade and development negotiations to develop an international commodity policy. An extensive summary of work on commodity futures trading is provided by JAMES B. WOY, *Commodity Futures Trading: A Bibliographic Guide* (1976).

Securities trading. *The Spicer & Oppenheim Guide to Securities Markets Around the World* (1988); and PAUL STONHAM, *Major Stock Markets of Europe* (1982), are good general surveys of world stock exchanges. ROBERT SOBEL, N.Y.S.E.: *A History of the New York Stock Exchange, 1935-1975* (1975), is a readable survey of this important exchange. See also JOEL SELIGMAN, *The Transformation of Wall Street: A History of the Securities and Exchange Commission and Modern Corporate Finance* (1982). Reference manuals include RICHARD J. TEWELES and EDWARD S. BRADLEY, *The Stock Market*, 5th ed. (1987); and FRANK G. ZARB and GABRIEL T. KERESKES (eds.), *The Stock Market Handbook* (1970). WILLIAM J. BAUMOL, *The Stock Market and Economic Efficiency* (1965), is an interesting effort to apply economic theory to the securities market. ARTHUR STONE DEWING, *The Financial Policy of Corporations*, 2 vol., 5th ed. (1953), is a classic treatise on financial policy, particularly useful for historical and statistical purposes. HUGH BULLOCK, *The Story of Investment Companies* (1959), recounts the development of mutual funds. VINCENT P. CAROSSO, *Investment Banking in America* (1970), provides a relatively thorough history. JOHN W. HAZARD and MILTON CHRISTIE, *The Investment Business* (1964), readably condenses the landmark

U.S. Securities and Exchange Commission's special study of the securities markets. Investment texts include HARRY C. SAUVAIN, *Investment Management*, 4th ed. (1973); and JEROME B. COHEN, EDWARD D. ZINBARG, and ARTHUR ZEIKEL, *Investment Analysis and Portfolio Management*, 5th ed. (1987), a comprehensive work. *Graham and Dodd's Security Analysis*, 5th ed. by SIDNEY COTTLE, ROGER F. MURRAY, and FRANK E. BLOCK (1988), is a classic work that led to the development of the field of security analysis. RICHARD W. JENNINGS, HAROLD MARSH, JR., and JOHN C. COFFEE, JR. (eds.), *Securities Regulation*, 7th ed. (1992), is a leading textbook dealing with the legal background of securities regulation. LOUIS LOSS and JOEL SELIGMAN, *Securities Regulation*, 3rd ed. (1989-), is a classic work kept up-to-date with supplements that delves into all aspects of the U.S. federal regulation of securities and securities markets.

Money market. The development and operation of the international capital market is addressed by M.S. MENDELSON, *Money on the Move: The Modern International Capital Market* (1980). Reference works include MARCIA STIGUM, *The Money Market*, 3rd ed. (1990), comprehensive and readable; and GUNTER DUFFEY and IAN H. GIDDY, *The International Money Market* (1978). TIMOTHY Q. COOK and TIMOTHY D. ROWE (eds.), *Instruments of the Money Market*, 6th ed. (1986), explains such key instruments as Eurodollars, treasury securities, and

federal funds. DAVID M. DARST, *The Handbook of the Bond and Money Markets* (1981), is a practical guide. Money markets in countries in Asia and the Pacific are studied by ARON VINER, *Inside Japanese Financial Markets* (1988); YOSHIO SUZUKI, *Money and Banking in Contemporary Japan*, trans. from the Japanese (1980), analyzing Japan's participation in international capital markets; *The Japanese Financial System* (1978), published by the Bank of Japan, a brief description of financial institutions, financial markets, and characteristics of the financial structure; and MICHAEL T. SKULLY (ed.), *Financial Institutions and Markets in the Far East* (1982), *Financial Institutions and Markets in Southeast Asia* (1984), *Financial Institutions and Markets in the Southwest Pacific* (1985), and *Financial Institutions and Markets in the South Pacific* (1987).

The markets and social welfare. The classic appraisal of the market from the standpoint of social welfare is A.C. PIGOU, *The Economics of Welfare*, 4th ed. (1962). Appraisals of welfare economics include I.M.D. LITTLE, *A Critique of Welfare Economics*, 2nd ed. (1957, reissued 1970); J. de V. GRAAFF, *Theoretical Welfare Economics* (1957, reissued 1975); and MAURICE DOBB, *On Economic Theory and Socialism* (1955, reissued 1972). THORSTEIN VEBLEN, *The Place of Science in Modern Civilization, and Other Essays* (1919, reprinted 1990), most directly expresses his critique of the market ideology. (Ed.)

Marseille

Founded more than 2,500 years ago, the port city of Marseille (English conventional spelling Marseilles) has a history of vigorous independence asserted against central authority in a variety of forms. It retained its status as a free city even after falling to Julius Caesar's troops in the 1st century BC, and after centuries of decline it was revived and allowed great independence under the local control of the viscounts of Provence in the 10th–14th centuries. After Provence joined the Kingdom of France in the 15th century, Marseille retained a separate administration and continually engaged in spirited revolt against kings or governments that threatened its liberties. It was for this reason that in 1800, when France was divided into the present administrative *départements*, Marseille was only reluctantly granted its status as capital of the Bouches-du-Rhône.

Frenchmen elsewhere, convinced that the Mediterranean climes of Provence could never be fully integrated into either the French realm or the Gallic spirit, long looked upon Marseille as a sort of folkloric institution: a place of comic anecdote and dialect, with a seasoning of picturesque criminality; a place where the citizens played a peculiar form of outdoor bowling known as *pétanque*, concocted the glorious garlic- and saffron-flavoured fish stew

known as bouillabaisse, and consumed rich, savory, absinthe-like Provençal pastis.

By whatever proportion fact may have been coloured with myth in its image, Marseille undoubtedly forms a major element in the economic and social structure of France. With Aix-en-Provence it forms the second largest urban agglomeration in France, and in association with the outport of Fos-sur-Mer, about 23 miles (37 kilometres) to the northwest, it is the country's largest seaport. Under the Socialist mayor Gaston Defferre, whose administration, from 1953 until his death in 1986, was the longest in its history, Marseille experienced major transformation—a process that is still continuing.

Marseille remains the capital of the Bouches-du-Rhône *département* and is also the administrative and commercial capital of Provence-Alpes-Côte d'Azur, one of France's fastest growing *régions*. It is situated on the Mediterranean's Gulf of Lion within a semicircle of limestone hills and lies 536 miles south-southeast of Paris by rail and 218 miles southeast of Lyon. The city proper has an area of about 93 square miles (241 square kilometres), and the metropolitan area covers some 365 square miles (946 square kilometres).

This article is divided into the following sections:

Physical and human geography 526

The landscape 526

The city site

Climate

The city layout

The people 528

The economy 528

Industry

Commerce and finance

Transportation

Administration and social conditions 529

Government

Services

Health

Education

Cultural life 529

History 529

The early period 529

Antiquity and the Middle Ages

Uneasy union with France

Era of expansion

The modern city 530

Bibliography 530

Physical and human geography

The character of Marseille has been determined to a great extent by geographic location. Its natural harbour, sheltered by a semicircle of hills and close to the estuary of the Rhône River, offered its first settlers the prospect of linking the Mediterranean seaways with northern Europe across a land that, in classical times, was made largely impassable by forests. The trading port founded by Greeks from the city of Phocaea in about 600 BC was to attract both settlers and visitors. The first of these account for a heterogeneous population and the second for services designed to cater to seamen and merchants. Marseille has the oldest chamber of commerce in France, established in 1599. It is a city of mosques and synagogues, besides many varieties of Christian churches. Its bars and brothels have been a magnet for dishonest dealings, and its waterfront still evokes the romance of a gateway to distant lands.

The city's most enduring characteristic, however, has been its readiness to welcome change. Its architecture preserves few vestiges of the past. Some landmarks, such as the transporter bridge that crossed the Old Port and the Panier district north of the harbour, were destroyed by the German occupation forces in 1943 and 1944. But more change has been wrought by the Marseillais themselves. Despite the legend that attaches to them, they are an unsentimental people open to new ideas.

For centuries, Marseille's mixed population and its inclination to political dissidence had made the city seem both foreign and marginal to French life and culture. After World War II, however, it was able to develop as a major European port and industrial centre. The city that had

been the starting point for past colonial enterprises bore a large share of the aftermath of French colonialism. It proved remarkably successful in absorbing new waves of immigrants, notably the former European colonists who crossed the Mediterranean from North Africa after Algeria won its independence in 1962. The building of the industrial complex at Fos-sur-Mer also attracted many thousands of migrant workers from North Africa in the 1960s, posing further problems of housing the city's immigrant population. Since this period, demographic growth has slowed. However, successive economic crises since the 1970s, combined with major restructuring by large industrial groups, have created a persistent unemployment, particularly among immigrants. This situation has been aggravated by racial hostility toward immigrants and by their marked concentration in particular districts of the city, especially certain northern suburbs. In an attempt to resolve such problems, these areas have become priority targets for the government's urban rehabilitation programs.

THE LANDSCAPE

The city site. Marseille lies in a sheltered depression surrounded by hills, which have inhibited the development of suburbs. The Old Port (Vieux-Port) is a natural harbour and one of the most westerly of the inlets along the rocky coastline characteristic of the northeastern Mediterranean; farther west, beyond the large tidal lake called the Berre Lagoon (Étang de Berre), the shoreline flattens out. There the sandy dunes of the Gulf of Fos and the Camargue region in the Rhône's delta were less attractive to early mariners and were only later seen as offering possibilities for development.

The
Old Port

Marseille's natural port was extended in Roman times and again from the 16th century onward to accommodate increased traffic and larger ships. By the 19th century the Old Port was insufficient. An artificial basin at La Joliette, built on the bay just outside of the Old Port, began operation in the mid-1840s, and five additional basins were subsequently built along a five-mile stretch of the bay. Further expansion was also undertaken to the west of the city, with the creation in 1863 of Port-Saint-Louis-du-Rhône. Eventually, in 1965, work began on the development of the port complex at Fos-sur-Mer. The port opened in 1968, and work on the accompanying industrial estate continued into the 1970s.

Marseille's hinterland consists of a chain of mountains, known as the Étoile Chain, which leads northward toward Aix-en-Provence (formerly Marseille's rival as capital of the region) and to Mount Sainte-Victoire. The slopes around Aix are devoted to vineyards, which produce the wines of the Côtes de Provence ("Hills of Provence"). The Étoile Chain has put a limit on the northward expansion of the city, with the result that development has "leapfrogged" this natural barrier in favour of the eastern shores of the Berre Lagoon, around the suburbs of Mari-gnane and Vitrolles. The eastern boundaries of the city also have been extended outward by the pressures of urban growth, particularly along the line of the Huveaune valley.

Climate. The climate of Marseille is not typical of the Mediterranean region as a whole. The lowest rainfall and highest temperatures are found in the hot, dry months of summer, but rainfall reaches a peak in spring and autumn, rather than in winter. The coldest months are December and January, when there is some frost, but otherwise the winter is mild. During the summer months the temperature rises to levels that would be unpleasant were it not for the sea breeze. In winter, Marseille is particularly subject to the dry, cold northwest wind known as the Mistral, which blows down the Rhône Valley, at times with considerable force.

The city layout. The popular area of Marseille was the seedy district, north of the Old Port, known as the Panier, which was destroyed in 1943. The more prosperous middle-class districts developed in the 19th century to the south around the rue Paradis and the avenue du Prado.

The period following World War II saw various schemes to develop the city, including the Unité d'Habitation, an 18-story residential block that expressed the architect Le Corbusier's ideal of urban family lodging. The block was intended, when completed in 1952, to be one of six such units; it is now surrounded by luxury apartment buildings. Less attractive are the high-rise, working-class housing blocks developed after 1960, which have come to house a large number of immigrant families. Since the 1970s the municipal authorities have undertaken a massive program of rebuilding in various parts of the city, restoring some of the Panier and turning the district around the Old Port into a largely pedestrian area. The area to the north of the Old Port, extending eastward toward the main railway station (Gare Saint-Charles), is the focus of a major program of urban renewal known as the Euroméditerranée, designed to refurbish existing buildings, provide new housing and office floor space, and create a new university site.

From the historic centre of Marseille at the Old Port, the thoroughfare of La Canebière climbs eastward up the hill; its name is a corruption of a Latin word for hemp, recalling Marseille's importance as a source of hemp and supplier of hemp rope in the Middle Ages. Thronged by people from around the world, La Canebière is the best-known commercial street in Marseille. Its starting point is marked by one of the most imposing public buildings in the city, the Bourse, which houses the Chamber of Commerce and a maritime museum.

The Bourse

Behind the Bourse, building operations in 1967 for a new retail and office complex uncovered a section of the Hellenistic ramparts of Massalia. Excavated by archaeologists, the site, dating from the 3rd and 2nd centuries BC, was found to consist of walls and towers and three sections of Roman road. The ancient port was also excavated. Nearby, close to City Hall on the edge of the port, is the Museum of Roman Warehouses (Musée des Docks Romains), which displays storage jars and other remains of commerce under Roman domination and traces the subsequent history of the port.

The port entrance is guarded by the Fort Saint-Jean, a 13th-century command post of the Knights Hospitaller of St. John of Jerusalem; some ruins remain, along with a tower built in the mid-15th century by King René of

By courtesy of the French Government Tourist Office



The Old Port and the sanctuary of Notre-Dame-de-la-Garde (centre), Marseille.

Provence. The extant fortress, dating from the 17th century, was part of a nationwide system of defenses. The other side of the harbour entrance is occupied by Fort Saint-Nicolas. In the harbour itself lie the Frioul Islands, on which the city has developed a large centre for water sports. Between these islands and the mainland is the Château d'If, a fortified island where Alexandre Dumas's fictional count of Monte Cristo and large numbers of all-too-real political prisoners were incarcerated.

Other historic buildings are located around the Old Port. In the Place de la Major, the old cathedral of la Major, built on the ruins of a temple of Diana, dates from the 11th century; it was partially dismantled to make way for the eight-domed structure that in 1852 replaced it as the city's cathedral. The dome and supporting arches of the old cathedral are fine examples of Provençal Romanesque stonemasonry.

Old
Charity
Hospital

Nearby is the Old Charity Hospital (Hospice de la Vieille Charité), built between 1660 and 1750. The interior courtyard surrounds a chapel by Pierre Puget, regarded as the most powerful of French Baroque sculptors. Close by is the Hôtel Dieu, the oldest hospital in the city, built at the end of the 16th century. The principal building, by Jules Hardouin-Mansart, was erected 200 years later and still serves its original function. Almost next door, the bell tower of the vanished church of Accoules, a 14th-century spire mounted on a 12th-century tower, marks the centre of Old Marseille.

On the opposite side of the port stands the crenellated, square-towered basilica of Saint-Victor, dating from the 11th to the 14th century; it once was attached to an abbey founded about 413 by St. John Cassian to commemorate a 3rd-century martyr and patron saint of sailors and millers. When Saint-Victor was built, the abbey was a temporal power of considerable extent, ruling properties in Spain, Sardinia, and the hinterlands of France.

High on the hill over the south side of the Old Port stands the celebrated Notre-Dame-de-la-Garde, a sanctuary honoured from the 8th century. Its present structure was built in 1853-64; its steeple, crowned by a 30-foot (nine-metre) gilded statue of the Virgin, rises 150 feet over the hillside.

THE PEOPLE

Marseille's population, drawn from all parts of the Mediterranean and from elsewhere in Europe and Africa, has always been mixed, so that it has never been possible to talk of a "typical" Marseillais. In 1880, for example, more than one in six of the inhabitants of the city was foreign. New residents have created a diverse pattern, sometimes concentrated in certain districts, such as the Muslim quarter that grew up during the 1970s north of La Canebière, and sometimes specializing in particular trades or professions. Certain groups—Jews, Greeks, Armenians—have their own community leaderships, which have semiofficial recognition.

Former colonials have had a strong impact on the community, and Marseille has always attracted Corsicans (including the Bonaparte family during the French Revolution). Manual labour is increasingly performed by North Africans or Africans who arrive from former colonies. There are marked social contrasts within the city. La Canebière forms an approximate dividing line between the working-class, often run-down areas of the north and the more affluent and salubrious districts of the south.

THE ECONOMY

Industry. Marseille itself has never been a major industrial centre; historically, its importance has been much more in trade and commerce. Nevertheless, certain industries did develop in Marseille. The oldest, founded in the 15th century, was the manufacture of soap from olive oil produced in the surrounding district. Other activities included food processing (linked to both imported agricultural products and those originating from the surrounding region), shipbuilding and ship repair, metallurgy, clothing, chemicals, and precision engineering. Many of these industries have either disappeared (as in the case of shipbuilding) or been reduced in importance through loss of markets or transfer to the city's periphery. Heavy industry

(oil refining and petrochemicals) grew up around the Berre Lagoon in the 1950s following the building of an outpost at Lavéra capable of receiving large oil tankers. This trend was accelerated from the late 1960s with the opening of the Fos port-industrial complex and with the addition of more petrochemical plants as well as steelworks. The majority of these installations use raw materials that enter through the port of Fos, and some of their finished products also leave by sea. The industrial zone is also directly linked to the national rail and highway networks, to the South European Pipeline, and to the Rhône inland waterway.

Lighter industrial development, warehousing, and transport-related activities have also greatly expanded north to the outlying districts of Marignane (site of Marseille's international airport) and Vitrolles. A similar trend is evident to the east along the Huveaune valley in the direction of Aubagne. Within Marseille itself a number of new industrial and related service activities have become established in fields such as electronics, data processing, telecommunications, and biomedicine. New sites have also been developed, including the Château-Gombert science park in the city's northeastern suburbs. The city's maritime location and traditions have also led to the growth of industries and services in offshore exploration and engineering.

The port complex of Marseille-Fos is the largest in France and among the largest in Europe. It is administered by the Port Autonome de Marseille ("Autonomous Port of Marseille"), a financially autonomous state enterprise that is responsible for the construction, administration, and maintenance of the industrial zones at Fos and Lavéra and the port facilities at Marseille, Lavéra, Caronte, Fos, and Port-Saint-Louis-du-Rhône. In addition to administering Marseille-Fos, the Port Autonome de Marseille also provides advice, information, and planning services for port authorities around the world. It is financed by rents, taxes, and fees for services, with state aid for investment in construction of harbours and quays.

In recent years the commercial traffic of the port complex has exceeded 90 million tons annually. The majority of this total is imports, mostly crude oil. Other imports include refined oil products, liquified natural gas, chemical products, and raw materials for the steel and aluminum industries. Exports consist mostly of refined oil products, chemicals, and steel. Containerized traffic of general merchandise is rising.

The different port zones have become increasingly specialized. Marseille itself handles roll-on/roll-off traffic (of both passengers and freight, principally to Corsica and North Africa), visits of cruise liners, and some bulk food products. Lavéra specializes in petroleum and chemical products, and Fos handles oil, other dry and liquid bulk cargoes, and containers. At Marseille there are also large ship repair yards, though their importance has greatly diminished.

Commerce and finance. Because of its geographic position and its commercial importance, Marseille has long been able to attract foreign capital. In the mid-19th century two major banks were established: the Société de Crédit Foncier de Marseille in 1852 and the Société Marseillaise de Crédit Industriel et Commercial et de Dépôts in 1865. Local money was directed both into industrial projects in the region and into foreign ventures such as the Suez Canal Company. But the disintegration of the French empire during the 20th century helped to accelerate the decline of Marseille as a financial centre. In this respect the relative lack of regional headquarters of large multinational firms represents another weakness in the city's economy.

Nonetheless, as capital of the Provence-Alpes-Côte d'Azur *région*, Marseille has considerable financial influence in the public sector. The Centre Méditerranéen de Commerce International (Mediterranean Centre for International Trade) was opened in 1983 to reassert the position of Marseille as a commercial and financial centre. It is anticipated that the Euroméditerranée complex will further enhance the city's role as a business capital.

Transportation. Marseille has good external connections. Two highways provide access to the north, and another highway reaches the city from the east. The high-

Port
develop-
ment

Innovations in transportation facilities

speed train (TGV; Train à Grande Vitesse), running on purpose-built track, makes it possible to reach Lyon in one hour and Paris in three. To the north of the city, the Marseille-Provence Airport (France's third-ranking airport for passenger traffic, after Paris and Nice) provides links to several destinations in France, Europe, and North Africa.

Within the city, movement is more problematic. Road congestion is severe, despite a tunnel under the centre linking the northern and eastern highways. Public transportation has been improved with the introduction of two underground metro lines and a surface tramway serving part of the eastern suburbs.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. The city government consists of a popularly elected municipal council. The council keeps very much alive the historical tradition of local independence in spite of the intimate involvement of many national ministries in the financing and planning of projects throughout the area.

The city is divided into 16 *arrondissements*, but for the purposes of local government these are grouped into eight *secteurs*, which elect mayors. In addition to the eight city halls, one for each *secteur*, there are two "mini city halls" in each *arrondissement*. The city mayor is assisted by a local government of 27 *adjoints*, each with responsibility for a particular facet of government, such as town planning, culture, finance, employment, or transport, and by delegate councillors who assist the *adjoints* or undertake more detailed responsibilities.

Services. The city's *adjoints* oversee the main urban services administered by the local authorities: lighting, refuse disposal, relations with the police and fire services, and so on. An unemployment rate above the national average and a large population of immigrant workers has exacerbated the problem of providing public services. In summer the region is particularly threatened by forest fires, and Marseille is the centre from which fire fighting is coordinated. A fleet of specially equipped airplanes is stationed at the airport.

Health. There are two teaching hospital complexes in Marseille, the North Hospital and the Timone Hospital. The centre for the study of tropical medicine at the Michel Lévy Hospital is well known. A computer centre in the suburb of Luminy links the region's hospitals.

Education. Three universities have sites in the city. The University of Aix-Marseille I offers courses in the sciences in Marseille (with courses in the arts and social sciences offered in Aix-en-Provence). The University of Aix-Marseille II has its faculty of medicine in Marseille, and the University of Aix-Marseille III has units in the sciences and engineering in the city. Both Aix and Marseille also have technical universities. In addition, there is a series of graduate schools specializing in fields such as physics, management, and engineering.

CULTURAL LIFE

Marseille has several museums, including a very popular Children's Museum. The Museum of Old Marseille was installed in 1960 next to the City Hall in Diamond House (La Maison Diamantée), so called because of its 16th-century facade of projecting diamond-shaped stone lozenges. The Cantini Museum, close to the rue Paradis, east of the Old Port, has a fine collection of Oriental art, of local pottery, and of modern paintings and sculptures.

Marseille has a number of historic sites and monuments, including its most famous landmark, the basilica of Notre-Dame-de-la-Garde, perched high above the city. It has an opera house, is an important centre for the theatre and music, and has a national dance school with a national ballet company. All of these activities are administered by the Municipal Office of Culture. The city has tried to reflect the diverse cultures of Marseille by encouraging exchanges with artists and companies from Algeria and other North African countries. During the 1930s, Marcel Pagnol founded a film studio in Marseille that made the city, for a time, France's only centre of the industry outside the Paris region.

As building has increased within Marseille, more attention has been paid to the conservation and development of

the municipality's parks and playgrounds. There are a large number of municipal sports centres and swimming pools, an outdoor theatre, and public beaches (centred around the Prado district and its aquarium). The parks of the Château du Pharo, Château Borély, and Palais Longchamp are extensive.

History

THE EARLY PERIOD

The oldest of the large French cities, Marseille was founded as Massalia (Massilia) by Greek mariners from Phocaea in Asia Minor about 600 BC. Archaeological finds exhibited in the Museum of Antiquities in the 18th-century Château Borély suggest that Phoenicians had settled there even earlier.

Antiquity and the Middle Ages. The Massalians spread trading posts inland as well as along the coasts, westward to Spain, and eastward to Monaco, founding the present cities of Arles, Nice, Antibes, Agde, and La Ciotat. Their coins have been found across France and through the Alps as far as the Tirol. In the 4th century BC a Massalian, Pytheas, visited the coasts of Gaul, Britain, and Germany, and a Euthymenes is said to have navigated the west coast of Africa as far south as Senegal.

When their great trade rivals, the Carthaginians, fought the Romans in the Punic Wars, Marseille supported the Romans and received help in subduing the native tribes of Liguria. When Pompey and Julius Caesar clashed, Marseille took Pompey's side and subsequently fell to Caesar's lieutenant Trebonius in 49 BC. Although stripped of dependencies, it was permitted to retain its status as a free city in recognition of past services. For some time the city remained the last centre of Greek learning in the West, but, eventually, it declined almost to extinction. After centuries of invasion and epidemic, it became little more than a huddle of nearly abandoned ruins.

In the 10th century, under the protection of successive viscounts of Provence, the area was repopulated, and it found new prosperity as a shipping and staging point for the Crusades. Gradually, the town bought up the rights of the viscounts, and, at the beginning of the 13th century, it formed a republic around the Old Port, though the upper part of the city and its southern suburb remained under ecclesiastical jurisdiction. The counts of Provence allowed the city great independence, and only in 1245 and 1256 did Charles of Anjou force acknowledgement of his sovereignty. After Marseille was sacked by Alfonso V of Aragon in 1423, King René of Provence, whose winter residence was there, restored prosperity to the city.

Uneasy union with France. When Provence, including Marseille, became part of the kingdom of France in 1481, the city preserved a separate administration directed by royal officials. During the 16th-century wars of religion, Marseille was fanatically Roman Catholic and long refused to recognize Henry IV as king because, until his conversion to Catholicism and accession to the French throne, he had been leader of the Protestants. During the Fronde, a movement in 1648-53 that opposed royal absolutism, the city sought to conserve its ancient liberties and rose against Louis XIV, who in 1660 came in person, breached the walls, and subdued the revolt. To discourage further manifestations of independence, the king planted Fort Saint-Nicolas at the southern extremity of the Old Port. In the same year, the city pushed inland to the west beyond its walls. A few buildings constructed in this expansion still survive in the area around the Cours Belsunce—a major thoroughfare—and the Préfecture.

Marseille joined enthusiastically in the French Revolution. Some 500 volunteers marching to Paris in 1792 sang "The War Song of the Rhine Army," which had been composed in Strasbourg in the late 18th century. The song, which thrilled the crowds along the route of march, was renamed "La Marseillaise" and became the national anthem of France. As the early federalist concepts were washed away in the blood of the Reign of Terror, however, the city revolted against the ruling Convention. Quickly mastered by force of arms, it was officially designated as "the city without a name." When its commerce was al-

Decline and resuscitation

Marseille and the Revolution

most destroyed by the maritime blockade of the Continent directed against Napoleon, Marseille became bitterly anti-Bonapartist and hailed the Bourbon restoration. Under Napoleon III, however, it remained stubbornly republican.

Era of expansion. During the second half of the 19th century, Marseille was expanded as the "port of empire," gaining impetus from the elimination of the Barbary pirates (1815–35), the conquest of Algeria (1830), and the inauguration of the Suez Canal (1869). The great avenues and many of the monuments of the city were constructed in this period. A serious water problem was solved by a project (1837–48) that brought water from the Durance River. The distribution reservoir above the city was disguised as the Longchamp Palace, containing the Museum of Natural History and the Museum of Fine Arts. The Château du Pharo, one of the city's principal landmarks, was built as a villa for Napoleon III and Eugénie at the edge of the bay beyond the Old Port, but it was never occupied by the imperial couple. The Bourse, an imposing colonnaded structure on La Canebière, was built in 1852–60 to house the Chamber of Commerce.

THE MODERN CITY

The city, occupied by the German army from November 1942 until August 1944, continued to be an active centre of the French Resistance movement, which partly explains the German decision to dynamite the Panier district and the Old Port in 1943. Further destruction was caused by German mines in August 1944.

The postwar history of Marseille is initially the story of the rebuilding of areas damaged in the war and the development of industrial and port complexes around Fos. It is also the story of a redistribution of population and economic activity. Central districts have lost population and have experienced industrial decline, particularly in areas

adjacent to the 19th-century port, while peripheral zones have expanded massively, acquiring large housing estates and new industrial and commercial parks. Marseille remains an important economic centre, though its influence in southeastern France is rivaled by the city of Lyon.

BIBLIOGRAPHY

General works: Chapters on Marseille, including maps, photographs, and site descriptions, can be found in numerous tour guides, such as *The Green Guide: Provence* (2000), by Michelin Travel Publications; and *Provence and the Côte d'Azur*, updated ed. (2001), in the Fodor's Guides series. ARCHIBALD LYALL, *The Companion Guide to the South of France*, 3rd ed. rev. and expanded by A.N. BRANGHAM (1972, reprinted 1996); and IAN B. THOMPSON, *The Lower Rhone and Marseille* (1975), include discussions of the topography, principal monuments, and cultural institutions of Marseille. M.F.K. FISHER, *A Considerable Town* (1978), presents vignettes on the city and its people by a noted American food and travel writer. ANDRÉ REMACLE, *Marseille à coeur ouvert* (1981), gives the personal view of a Marseille journalist and deals with the character and history of the city.

PAUL CARRÈRE and RAYMOND DUGRAND, *La Région Méditerranéenne*, 2nd rev. ed. (1967), deals with the city's industrial region. GILBERT ROCHU, *Marseille, les années Defferre* (1983), is a critical analysis of the administration of the city in the period since World War II. More recent works on the development of Marseille include A. MÉDAM, *Blues Marseille* (1995); DOMINIQUE BECQUART, *Marseille, 25 ans de planification urbaine* (1994); JEAN-LUCIEN BONILLO et al., *Marseille, ville & port* (1992); JEAN VIARD, *Marseille, une ville impossible* (1995); and MARCEL RONCAYOLO, *Marseille: les territoires du temps* (1996).

History: Two comprehensive histories in French are GASTON RAMBERT (ed.), *Histoire du commerce de Marseille*, 6 vol. (1949–59); and RAOUL BUSQUET, *Histoire de Marseille*, new ed., rev. and updated by CONSTANT VAUTRAVERS (1998). LUCIEN GAILLARD, *Marseille sous l'occupation* (1982), covers the occupation years of World War II and is illustrated with contemporary photographs and reproductions of documents.

(B.E./R.C.Bu./J.N.T.)

Marx and Marxism

Karl Marx, revolutionary, sociologist, historian, and economist, was the author (with Friedrich Engels) of *Manifest der Kommunistischen Partei* (1848), commonly known as *The Communist Manifesto*, the most celebrated pamphlet in the history of the socialist movement, as well as of its most important book, *Das Kapital*. These writings and others by Marx and Engels form the basis of the body of thought and belief known as Marxism.

This article deals with Marx's life, his thinking, his accomplishments, and the development of Marxist theory. See the article SOCIO-ECONOMIC DOCTRINES AND REFORM MOVEMENTS, MODERN, for a full treatment of socialism and communism. For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, section 541, and the *Index*.

The article is divided into the following sections:

Life and works of Marx	531	The work of Kautsky and Bernstein	
Early years	531	The radicals	
Brussels period	532	The Austrians	
Early years in London	533	Russian and Soviet Marxism	539
Role in the First International	533	Lenin	
Last years	534	The dictatorship of the proletariat	
Character and significance	534	Stalin	
Marxism	535	Trotskyism	
The thought of Karl Marx	535	Variants of Marxism	541
Historical materialism		Maoism	
Analysis of society		Marxism in Cuba	
Analysis of the economy		Marxism in the Third World	
Class struggle		Marxism in the West	
The contributions of Engels		Major works	542
German Marxism after Engels	538	Bibliography	542

Life and works of Marx

EARLY YEARS

Karl Heinrich Marx was born on May 5, 1818, in the city of Trier in the Rhine province of Prussia, now in Germany. He was the oldest surviving boy of nine children. His father, Heinrich, a successful lawyer, was a man of the Enlightenment, devoted to Kant and Voltaire, who took part in agitations for a constitution in Prussia. His mother, born Henrietta Pressburg, was from Holland. Both parents were Jewish and were descended from a long line of rabbis, but, a year or so before Karl was born, his father—probably because his professional career required it—was baptized in the Evangelical Established Church. Karl was baptized when he was six years old. Although as a youth Karl was influenced less by religion than by the critical, sometimes radical social policies of the Enlightenment, his Jewish background exposed him to prejudice and discrimination that may have led him to question the role of religion in society and contributed to his desire for social change.

Early
education

Marx was educated from 1830 to 1835 at the high school in Trier. Suspected of harbouring liberal teachers and pupils, the school was under police surveillance. Marx's writings during this period exhibited a spirit of Christian devotion and a longing for self-sacrifice on behalf of humanity. In October 1835 he matriculated at the University of Bonn. The courses he attended were exclusively in the humanities, in such subjects as Greek and Roman mythology and the history of art. He participated in customary student activities, fought a duel, and spent a day in jail for being drunk and disorderly. He presided at the Tavern Club, which was at odds with the more aristocratic student associations, and joined a poets' club that included some political activists. A politically rebellious student culture was, indeed, part of life at Bonn. Many students had been arrested; some were still being expelled in Marx's time, particularly as a result of an effort by students to disrupt a session of the Federal Diet at Frankfurt. Marx, however, left Bonn after a year and in October 1836 enrolled at the University of Berlin to study law and philosophy.

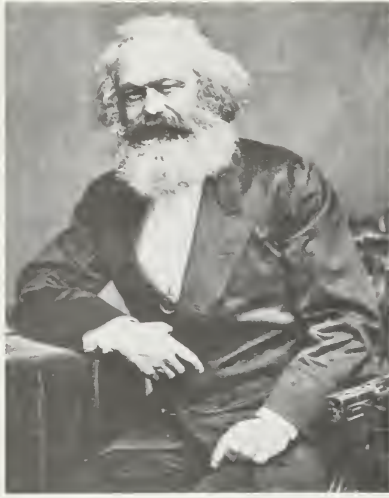
Marx's crucial experience at Berlin was his introduction to Hegel's philosophy, regnant there, and his adherence to the Young Hegelians. At first he felt a repugnance toward Hegel's doctrines; when Marx fell sick it was partially, as

he wrote his father, "from intense vexation at having to make an idol of a view I detested." The Hegelian pressure in the revolutionary student culture was powerful, however, and Marx joined a society called the Doctor Club, whose members were intensely involved in the new literary and philosophical movement. Their chief figure was Bruno Bauer, a young lecturer in theology, who was developing the idea that the Christian Gospels were a record not of history but of human fantasies arising from emotional needs and that Jesus had not been a historical person. Marx enrolled in a course of lectures given by Bauer on the prophet Isaiah. Bauer taught that a new social catastrophe "more tremendous" than that of the advent of Christianity was in the making. The Young Hegelians began moving rapidly toward atheism and also talked vaguely of political action.

The Prussian government, fearful of the subversion latent in the Young Hegelians, soon undertook to drive them from the universities. Bauer was dismissed from his post in 1839. Marx's "most intimate friend" of this period, Adolph Rutenberg, an older journalist who had served a prison sentence for his political radicalism, pressed for a deeper social involvement. By 1841 the Young Hegelians had become left republicans. Marx's studies, meanwhile, were lagging. Urged by his friends, he submitted a doctoral dissertation to the university at Jena, which was known to be lax in its academic requirements, and received his degree in April 1841. His thesis analyzed in a Hegelian fashion the difference between the natural philosophies of Democritus and Epicurus. More distinctively, it sounded a note of Promethean defiance:

Philosophy makes no secret of it. Prometheus' admission: "In sooth all gods I hate," is its own admission, its own motto against all gods. . . . Prometheus is the noblest saint and martyr in the calendar of philosophy.

In 1841 Marx, together with other Young Hegelians, was much influenced by the publication of *Das Wesen des Christentums* (1841; *The Essence of Christianity*) by Ludwig Feuerbach. Its author, to Marx's mind, successfully criticized Hegel, an idealist who believed that matter or existence was inferior to and dependent upon mind or spirit, from the opposite, or materialist, standpoint, showing how the "Absolute Spirit" was a projection of "the real man standing on the foundation of nature." Henceforth Marx's philosophical efforts were toward a combination of



Marx.

By courtesy of the trustees of the British Museum photograph J.R. Freeman & Co Ltd

Hegel's dialectic—the idea that all things are in a continual process of change resulting from the conflicts between their contradictory aspects—with Feuerbach's materialism, which placed material conditions above ideas.

In January 1842 Marx began contributing to a newspaper newly founded in Cologne, the *Rheinische Zeitung*. It was the liberal democratic organ of a group of young merchants, bankers, and industrialists; Cologne was the centre of the most industrially advanced section of Prussia. To this stage of Marx's life belongs an essay on the freedom of the press. Since he then took for granted the existence of absolute moral standards and universal principles of ethics, he condemned censorship as a moral evil that entailed spying into people's minds and hearts and assigned to weak and malevolent mortal powers that presupposed an omniscient mind. He believed that censorship could have only evil consequences.

Newspaper work

On Oct. 15, 1842, Marx became editor of the *Rheinische Zeitung*. As such, he was obliged to write editorials on a variety of social and economic issues, ranging from the housing of the Berlin poor and the theft by peasants of wood from the forests to the new phenomenon of communism. He found Hegelian idealism of little use in these matters. At the same time he was becoming estranged from his Hegelian friends for whom shocking the bourgeois was a sufficient mode of social activity. Marx, friendly at this time to the "liberal-minded practical men" who were "struggling step-by-step for freedom within constitutional limits," succeeded in trebling his newspaper's circulation and making it a leading journal in Prussia. Nevertheless, Prussian authorities suspended it for being too outspoken, and Marx agreed to coedit with the liberal Hegelian Arnold Ruge a new review, the *Deutsch-französische Jahrbücher* ("German-French Yearbooks"), which was to be published in Paris.

First, however, in June 1843 Marx, after an engagement of seven years, married Jenny von Westphalen. Jenny was an attractive, intelligent, and much-admired woman, four years older than Karl; she came of a family of military and administrative distinction. Her half-brother later became a highly reactionary Prussian minister of the interior. Her father, a follower of the French socialist Saint-Simon, was fond of Karl, though others in her family opposed the marriage. Marx's father also feared that Jenny was destined to become a sacrifice to the demon that possessed his son.

Four months after their marriage, the young couple moved to Paris, which was then the centre of socialist thought and of the more extreme sects that went under the name of communism. There, Marx first became a revolutionary and a communist and began to associate with communist societies of French and German workingmen. Their ideas were, in his view, "utterly crude and unintelligent," but their character moved him: "The brotherhood of man is no mere phrase with them, but a fact of life, and the nobility of man shines upon us from their work-

hardened bodies," he wrote in his so-called "Ökonomisch-philosophische Manuskripte aus dem Jahre 1844" (written in 1844; *Economic and Philosophic Manuscripts of 1844* [1959]). (These manuscripts were not published for some 100 years, but they are influential because they show the humanist background to Marx's later historical and economic theories.)

The "German-French Yearbooks" proved short-lived, but through their publication Marx befriended Friedrich Engels, a contributor who was to become his lifelong collaborator, and in their pages appeared Marx's article "Zur Kritik der Hegelschen Rechtsphilosophie" ("Toward the Critique of the Hegelian Philosophy of Right") with its oft-quoted assertion that religion is the "opium of the people." It was there, too, that he first raised the call for an "uprising of the proletariat" to realize the conceptions of philosophy. Once more, however, the Prussian government intervened against Marx. He was expelled from France and left for Brussels—followed by Engels—in February 1845. That year in Belgium he renounced his Prussian nationality.

BRUSSELS PERIOD

The next two years in Brussels saw the deepening of Marx's collaboration with Engels. Engels had seen at firsthand in Manchester, Eng., where a branch factory of his father's textile firm was located, all the depressing aspects of the Industrial Revolution. He had also been a Young Hegelian and had been converted to communism by Moses Hess, who was called the "communist rabbi." In England he associated with the followers of Robert Owen. Now he and Marx, finding that they shared the same views, combined their intellectual resources and published *Die heilige Familie* (1845; *The Holy Family*), a prolix criticism of the Hegelian idealism of the theologian Bruno Bauer. Their next work, *Die deutsche Ideologie* (written 1845–46, published 1932; *The German Ideology*), contained the fullest exposition of their important materialistic conception of history, which set out to show how, historically, societies had been structured to promote the interests of the economically dominant class. But it found no publisher and remained unknown during its authors' lifetimes.

During his Brussels years, Marx developed his views and, through confrontations with the chief leaders of the working-class movement, established his intellectual standing. In 1846 he publicly excoriated the German leader Wilhelm Weitling for his moralistic appeals. Marx insisted that the stage of bourgeois society could not be skipped over; the proletariat could not just leap into communism; the workers' movement required a scientific basis, not moralistic phrases. He also polemicized against the French socialist thinker Pierre-Joseph Proudhon in *Misère de la philosophie* (1847; *The Poverty of Philosophy*), a mordant attack on Proudhon's book subtitled *Philosophie de la misère* (1846; *The Philosophy of Poverty*). Proudhon wanted to unite the best features of such contraries as competition and monopoly; he hoped to save the good features in economic institutions while eliminating the bad. Marx, however, declared that no equilibrium was possible between the antagonisms in any given economic system. Social structures were transient historic forms determined by the productive forces: "The handmill gives you society with the feudal lord; the steammill, society with the industrial capitalist." Proudhon's mode of reasoning, Marx wrote, was typical of the petty bourgeois, who failed to see the underlying laws of history.

An unusual sequence of events led Marx and Engels to write their pamphlet *The Communist Manifesto*. In June 1847 a secret society, the League of the Just, composed mainly of emigrant German handicraftsmen, met in London and decided to formulate a political program. They sent a representative to Marx to ask him to join the league; Marx overcame his doubts and, with Engels, joined the organization, which thereupon changed its name to the Communist League and enacted a democratic constitution. Entrusted with the task of composing their program, Marx and Engels worked from the middle of December 1847 to the end of January 1848. The London Communists were already impatiently threatening Marx with

Collaboration with Engels

The Communist Manifesto

disciplinary action when he sent them the manuscript; they promptly adopted it as their manifesto. It enunciated the proposition that all history had hitherto been a history of class struggles, summarized in pithy form the materialist conception of history worked out in *The German Ideology*, and asserted that the forthcoming victory of the proletariat would put an end to class society forever. It mercilessly criticized all forms of socialism founded on philosophical "cobwebs" such as "alienation." It rejected the avenue of "social Utopias," small experiments in community, as deadening the class struggle and therefore as being "reactionary sects." It set forth 10 immediate measures as first steps toward communism, ranging from a progressive income tax and the abolition of inheritances to free education for all children. It closed with the words, "The proletarians have nothing to lose but their chains. They have a world to win. Workingmen of all countries, unite!"

Revolution suddenly erupted in Europe in the first months of 1848, in France, Italy, and Austria. Marx had been invited to Paris by a member of the provisional government just in time to avoid expulsion by the Belgian government. As the revolution gained in Austria and Germany, Marx returned to the Rhineland. In Cologne he advocated a policy of coalition between the working class and the democratic bourgeoisie, opposing for this reason the nomination of independent workers' candidates for the Frankfurt Assembly and arguing strenuously against the program for proletarian revolution advocated by the leaders of the Workers' Union. He concurred in Engels' judgment that *The Communist Manifesto* should be shelved and the Communist League disbanded. Marx pressed his policy through the pages of the *Neue Rheinische Zeitung*, newly founded in June 1849, urging a constitutional democracy and war with Russia. When the more revolutionary leader of the Workers' Union, Andreas Gottschalk, was arrested, Marx supplanted him and organized the first Rhineland Democratic Congress in August 1848. When the king of Prussia dissolved the Prussian Assembly in Berlin, Marx called for arms and men to help the resistance. Bourgeois liberals withdrew their support from Marx's newspaper, and he himself was indicted on several charges, including advocacy of the nonpayment of taxes. In his trial he defended himself with the argument that the crown was engaged in making an unlawful counterrevolution. The jury acquitted him unanimously and with thanks. Nevertheless, as the last hopeless fighting flared in Dresden and Baden, Marx was ordered banished as an alien on May 16, 1849. The final issue of his newspaper, printed in red, caused a great sensation.

EARLY YEARS IN LONDON

Expelled once more from Paris, Marx went to London in August 1849. It was to be his home for the rest of his life. Chagrined by the failure of his own tactics of collaboration with the liberal bourgeoisie, he rejoined the Communist League in London and for about a year advocated a bolder revolutionary policy. An "Address of the Central Committee to the Communist League," written with Engels in March 1850, urged that in future revolutionary situations they struggle to make the revolution "permanent" by avoiding subservience to the bourgeois party and by setting up "their own revolutionary workers' governments" alongside any new bourgeois one. Marx hoped that the economic crisis would shortly lead to a revival of the revolutionary movement; when this hope faded, he came into conflict once more with those whom he called "the alchemists of the revolution," such as August von Willich, a communist who proposed to hasten the advent of revolution by undertaking direct revolutionary ventures. Such persons, Marx wrote in September 1850, substitute "idealism for materialism" and regard

pure will as the motive power of revolution instead of actual conditions. While we say to the workers: "You have got to go through fifteen, twenty, fifty years of civil wars and national wars not merely in order to change your conditions but in order to change yourselves and become qualified for political power," you on the contrary tell them, "We must achieve power immediately."

The militant faction in turn ridiculed Marx for being a revolutionary who limited his activity to lectures on political economy to the Communist Workers' Educational Union. The upshot was that Marx gradually stopped attending meetings of the London Communists. In 1852 he devoted himself intensely to working for the defense of 11 communists arrested and tried in Cologne on charges of revolutionary conspiracy and wrote a pamphlet on their behalf. The same year he also published, in a German-American periodical, his essay "Der Achtzehnte Brumaire des Louis Napoleon" (*The Eighteenth Brumaire of Louis Bonaparte*), with its acute analysis of the formation of a bureaucratic absolutist state with the support of the peasant class. In other respects the next 12 years were, in Marx's words, years of "isolation" both for him and for Engels in his Manchester factory.

From 1850 to 1864 Marx lived in material misery and spiritual pain. His funds were gone, and except on one occasion he could not bring himself to seek paid employment. In March 1850 he and his wife and four small children were evicted and their belongings seized. Several of his children died—including a son Guido, "a sacrifice to bourgeois misery," and a daughter Franziska, for whom his wife rushed about frantically trying to borrow money for a coffin. For six years the family lived in two small rooms in Soho, often subsisting on bread and potatoes. The children learned to lie to the creditors: "Mr. Marx ain't upstairs." Once he had to escape them by fleeing to Manchester. His wife suffered breakdowns.

During all these years Engels loyally contributed to Marx's financial support. The sums were not large at first, for Engels was only a clerk in the firm of Ermen and Engels at Manchester. Later, however, in 1864, when he became a partner, his subventions were generous. Marx was proud of Engels' friendship and would tolerate no criticism of him. Bequests from the relatives of Marx's wife and from Marx's friend Wilhelm Wolff also helped to alleviate their economic distress.

Marx had one relatively steady source of earned income in the United States. On the invitation of Charles A. Dana, managing editor of *The New York Tribune*, he became in 1851 its European correspondent. The newspaper, edited by Horace Greeley, had sympathies for Fourierism, a Utopian socialist system developed by the French theorist Charles Fourier. From 1851 to 1862 Marx contributed close to 500 articles and editorials (Engels providing about a fourth of them). He ranged over the whole political universe, analyzing social movements and agitations from India and China to Britain and Spain.

In 1859 Marx published his first book on economic theory, *Zur Kritik der politischen Ökonomie* (*A Contribution to the Critique of Political Economy*). In its preface he again summarized his materialistic conception of history, his theory that the course of history is dependent on economic developments. At this time, however, Marx regarded his studies in economic and social history at the British Museum as his main task. He was busy producing the drafts of his magnum opus, which was to be published later as *Das Kapital*. Some of these drafts, including the *Outlines* and the *Theories of Surplus Value*, are important in their own right and were published after Marx's death.

ROLE IN THE FIRST INTERNATIONAL

Marx's political isolation ended in 1864 with the founding of the International Working Men's Association. Although he was neither its founder nor its head, he soon became its leading spirit. Its first public meeting, called by English trade union leaders and French workers' representatives, took place at St. Martin's Hall in London on Sept. 28, 1864. Marx, who had been invited through a French intermediary to attend as a representative of the German workers, sat silently on the platform. A committee was set up to produce a program and a constitution for the new organization. After various drafts had been submitted that were felt to be unsatisfactory, Marx, serving on a subcommittee, drew upon his immense journalistic experience. His "Address and the Provisional Rules of the International Working Men's Association," unlike his other writings, stressed the positive achievements of the cooperative

Participation in the events of 1848

Poverty in London

movement and of parliamentary legislation; the gradual conquest of political power would enable the British proletariat to extend these achievements on a national scale.

As a member of the organization's General Council, and corresponding secretary for Germany, Marx was henceforth assiduous in attendance at its meetings, which were sometimes held several times a week. For several years he showed a rare diplomatic tact in composing differences among various parties, factions, and tendencies. The International grew in prestige and membership, its numbers reaching perhaps 800,000 in 1869. It was successful in several interventions on behalf of European trade unions engaged in struggles with employers.

The Paris
Commune

In 1870, however, Marx was still unknown as a European political personality; it was the Paris Commune that made him into an international figure, "the best calumniated and most menaced man of London," as he wrote. When the Franco-German War broke out in 1870, Marx and Engels disagreed with followers in Germany who refused to vote in the Reichstag in favour of the war. The General Council declared that "on the German side the war was a war of defence." After the defeat of the French armies, however, they felt that the German terms amounted to agrandizement at the expense of the French people. When an insurrection broke out in Paris and the Paris Commune was proclaimed, Marx gave it his unswerving support. On May 30, 1871, after the Commune had been crushed, he hailed it in a famous address entitled *Civil War in France*:

History has no comparable example of such greatness. . . . Its martyrs are enshrined forever in the great heart of the working class.

In Engels' judgment, the Paris Commune was history's first example of the "dictatorship of the proletariat." Marx's name, as the leader of The First International and author of the notorious *Civil War*, became synonymous throughout Europe with the revolutionary spirit symbolized by the Paris Commune.

The advent of the Commune, however, exacerbated the antagonisms within the International Working Men's Association and thus brought about its downfall. English trade unionists such as George Odger, former president of the General Council, opposed Marx's support of the Paris Commune. The Reform Bill of 1867, which had enfranchised the British working class, had opened vast opportunities for political action by the trade unions. English labour leaders found they could make many practical advances by cooperating with the Liberal Party and, regarding Marx's rhetoric as an encumbrance, resented his charge that they had "sold themselves" to the Liberals.

The
struggle
with
Bakunin

A left opposition also developed under the leadership of the famed Russian revolutionary Mikhail Alexandrovich Bakunin. A veteran of tsarist prisons and Siberian exile, Bakunin could move men by his oratory, which one listener compared to "a raging storm with lightning, flashes and thunderclaps, and a roaring as of lions." Bakunin admired Marx's intellect but could hardly forget that Marx had published a report in 1848 charging him with being a Russian agent. He felt that Marx was a German authoritarian and an arrogant Jew who wanted to transform the General Council into a personal dictatorship over the workers. He strongly opposed several of Marx's theories, especially Marx's support of the centralized structure of the International, Marx's view that the proletariat class should act as a political party against prevailing parties but within the existing parliamentary system, and Marx's belief that the proletariat, after it had overthrown the bourgeois state, should establish its own regime. To Bakunin, the mission of the revolutionary was destruction; he looked to the Russian peasantry, with its propensities for violence and its uncurbed revolutionary instincts, rather than to the effete, civilized workers of the industrial countries. The students, he hoped, would be the officers of the revolution. He acquired followers, mostly young men, in Italy, Switzerland, and France, and he organized a secret society, the International Alliance of Social Democracy, which in 1869 challenged the hegemony of the General Council at the congress in Basel, Switz. Marx, however, had already succeeded in preventing its admission as an organized body into the International.

To the Bakuninists, the Paris Commune was a model of revolutionary direct action and a refutation of what they considered to be Marx's "authoritarian communism." Bakunin began organizing sections of the International for an attack on the alleged dictatorship of Marx and the General Council. Marx in reply publicized Bakunin's embroilment with an unscrupulous Russian student leader, Sergey Gennadiyevich Nechayev, who had practiced blackmail and murder.

Without a supporting right wing and with the anarchist left against him, Marx feared losing control of the International to Bakunin. He also wanted to return to his studies and to finish *Das Kapital*. At the congress of the International at The Hague in 1872, the only one he ever attended, Marx managed to defeat the Bakuninists. Then, to the consternation of the delegates, Engels moved that the seat of the General Council be transferred from London to New York City. The Bakuninists were expelled, but the International languished and was finally disbanded in Philadelphia in 1876.

Dissolution
of the
Inter-
national

LAST YEARS

During the next and last decade of his life, Marx's creative energies declined. He was beset by what he called "chronic mental depression," and his life turned inward toward his family. He was unable to complete any substantial work, though he still read widely and undertook to learn Russian. He became crotchety in his political opinions. When his own followers and those of the German revolutionary Ferdinand Lassalle, a rival who believed that socialist goals should be achieved through cooperation with the state, coalesced in 1875 to found the German Social Democratic Party, Marx wrote a caustic criticism of their program (the so-called Gotha Program), claiming that it made too many compromises with the status quo. The German leaders put his objections aside and tried to mollify him personally. Increasingly, he looked to a European war for the overthrow of Russian tsarism, the mainstay of reaction, hoping that this would revive the political energies of the working classes. He was moved by what he considered to be the selfless courage of the Russian terrorists who assassinated the tsar, Alexander II, in 1881; he felt this to be "a historically inevitable means of action."

Despite Marx's withdrawal from active politics, he still retained what Engels called his "peculiar influence" on the leaders of working-class and socialist movements. In 1879, when the French Socialist Workers' Federation was founded, its leader Jules Guesde went to London to consult with Marx, who dictated the preamble of its program and shaped much of its content. In 1881 Henry Mayers Hyndman in his *England for All* drew heavily on his conversations with Marx but angered him by being afraid to acknowledge him by name.

During his last years Marx spent much time at health resorts and even traveled to Algiers. He was broken by the death of his wife on Dec. 2, 1881, and of his eldest daughter, Jenny Longuet, on Jan. 11, 1883. He died in London, evidently of a lung abscess, on March 14, 1883.

Death

CHARACTER AND SIGNIFICANCE

At Marx's funeral in Highgate Cemetery, Engels declared that Marx had made two great discoveries, the law of development of human history and the law of motion of bourgeois society. But "Marx was before all else a revolutionist." He was "the best-hated and most-calumniated man of his time," yet he also died "beloved, revered and mourned by millions of revolutionary fellow-workers."

The contradictory emotions Marx engendered are reflected in the sometimes conflicting aspects of his character. Marx was a combination of the Promethean rebel and the rigorous intellectual. He gave most persons an impression of intellectual arrogance. A Russian writer, Pavel Annenkov, who observed Marx in debate in 1846 recalled that "he spoke only in the imperative, brooking no contradiction," and seemed to be "the personification of a democratic dictator such as might appear before one in moments of fantasy." But Marx obviously felt uneasy before mass audiences and avoided the atmosphere of factional controversies at congresses. He went to no

demonstrations, his wife remarked, and rarely spoke at public meetings. He kept away from the congresses of the International where the rival socialist groups debated important resolutions. He was a "small groups" man, most at home in the atmosphere of the General Council or on the staff of a newspaper, where his character could impress itself forcefully on a small body of coworkers. At the same time he avoided meeting distinguished scholars with whom he might have discussed questions of economics and sociology on a footing of intellectual equality. Despite his broad intellectual sweep, he was prey to obsessive ideas such as that the British foreign minister, Lord Palmerston, was an agent of the Russian government. He was determined not to let bourgeois society make "a money-making machine" out of him, yet he submitted to living on the largess of Engels and the bequests of relatives. He remained the eternal student in his personal habits and way of life, even to the point of joining two friends in a students' prank during which they systematically broke four or five streetlamps in a London street and then fled from the police. He was a great reader of novels, especially those of Sir Walter Scott and Balzac; and the family made a cult of Shakespeare. He was an affectionate father, saying that he admired Jesus for his love of children, but sacrificed the lives and health of his own. Of his seven children, three daughters grew to maturity. His favourite daughter, Eleanor, worried him with her nervous, brooding, emotional character and her desire to be an actress. Another shadow was cast on Marx's domestic life by the birth to their loyal servant, Helene Demuth, of an illegitimate son, Frederick; Engels as he was dying disclosed to Eleanor that Marx had been the father. Above all, Marx was a fighter, willing to sacrifice anything in the battle for his conception of a better society. He regarded struggle as the law of life and existence.

The influence of Marx's ideas has been enormous. Marx's masterpiece, *Das Kapital*, the "Bible of the working class," as it was officially described in a resolution of the International Working Men's Association, was published in 1867 in Berlin and received a second edition in 1873. Only the first volume was completed and published in Marx's lifetime. The second and third volumes, unfinished by Marx, were edited by Engels and published in 1885 and 1894. The economic categories he employed were those of the classical British economics of David Ricardo; but Marx used them in accordance with his dialectical method to argue that bourgeois society, like every social organism, must follow its inevitable path of development. Through the working of such immanent tendencies as the declining rate of profit, capitalism would die and be replaced by another, higher, society. The most memorable pages in *Das Kapital* are the descriptive passages, culled from Parliamentary Blue Books, on the misery of the English working class. Marx believed that this misery would increase, while at the same time the monopoly of capital would become a fetter upon production until finally "the knell of capitalist private property sounds. The expropriators are expropriated."

Marx never claimed to have discovered the existence of classes and class struggles in modern society. "Bourgeois" historians, he acknowledged, had described them long before he had. He did claim, however, to have proved that each phase in the development of production was associated with a corresponding class structure and that the struggle of classes led necessarily to the dictatorship of the proletariat, ushering in the advent of a classless society. Marx took up the very different versions of socialism current in the early 19th century and welded them together into a doctrine that continued to be the dominant version of socialism for half a century after his death. His emphasis on the influence of economic structure on historical development has proved to be of lasting significance.

Although Marx stressed economic issues in his writings, his major impact has been in the fields of sociology and history. Marx's most important contribution to sociological theory was his general mode of analysis, the "dialectical" model, which regards every social system as having within it immanent forces that give rise to "contradictions" (disequilibria) that can be resolved only by a new

social system. Neo-Marxists, who no longer accept the economic reasoning in *Das Kapital*, are still guided by this model in their approach to capitalist society. In this sense, Marx's mode of analysis, like those of Thomas Malthus, Herbert Spencer, or Vilfredo Pareto, has become one of the theoretical structures that are the heritage of the social scientist. (L.S.F./D.T.McL.)

Marxism

The term Marxism is used in a number of different ways. In its most essential meaning it refers to the thought of Karl Marx but is usually extended to include that of his friend and collaborator Friedrich Engels. There is also Marxism as it has been understood and practiced by the various socialist movements, particularly before 1914. Then there is Soviet Marxism as worked out by Lenin and modified by Stalin, which under the name of Marxism-Leninism became the doctrine of the communist parties set up after the Russian Revolution. Offshoots of this include Marxism as interpreted by the anti-Stalinist Leon Trotsky and his followers, Mao Zedong's (Mao Tse-tung's) Chinese variant of Marxism-Leninism, and various Third World Marxisms. There are also the post-World War II nondogmatic Marxisms that have modified Marx's thought with borrowings from modern philosophies, principally from those of Edmund Husserl and Martin Heidegger but also from Sigmund Freud and others.

THE THOUGHT OF KARL MARX

The written work of Marx cannot be reduced to a philosophy, much less to a philosophical system. The whole of his work is a radical critique of philosophy, especially of Hegel's idealist system and of the philosophies of the left and right post-Hegelians. It is not, however, a mere denial of those philosophies. Marx declared that philosophy must become reality. One could no longer be content with interpreting the world; one must be concerned with transforming it, which meant transforming both the world itself and men's consciousness of it. This, in turn, required a critique of experience together with a critique of ideas. In fact, Marx believed that all knowledge involved a critique of ideas. He was not an empiricist. Rather, his work teems with concepts (appropriation, alienation, praxis, creative labour, value, etc.) that he had inherited from earlier philosophers and economists, including Hegel, Johann Fichte, Kant, Adam Smith, David Ricardo, and John Stuart Mill. What uniquely characterizes the thought of Marx is that, instead of making abstract affirmations about a whole group of problems such as man, knowledge, matter, and nature, he examines each problem in its dynamic relation to the others and, above all, tries to relate them to historical, social, political, and economic realities.

Historical materialism. In 1859, in the preface to his *Contribution to the Critique of Political Economy*, Marx wrote that the hypothesis that had served him as the basis for his analysis of society could be briefly formulated as follows:

In the social production that men carry on, they enter into definite relations that are indispensable and independent of their will, relations of production which correspond to a definite stage of development of their material forces of production. The sum total of these relations of production constitutes the economic structure of society, the real foundation, on which rises a legal and political superstructure, and to which correspond definite forms of social consciousness. The mode of production in material life determines the general character of the social, political, and intellectual processes of life. It is not the consciousness of men which determines their existence; it is on the contrary their social existence which determines their consciousness.

Raised to the level of historical law, this hypothesis was subsequently called historical materialism. Marx applied it to capitalist society, both in *The Communist Manifesto* and *Das Kapital* and in other writings. Although Marx reflected upon his working hypothesis for many years, he did not formulate it in a very exact manner: different expressions served him for identical realities. If one takes the text literally, social reality is structured in the following way:

Economic
foundation
of society

*Das
Kapital*

The
dialectical
model

1. Underlying everything as the real basis of society is the economic structure (what in late 20th-century language is sometimes called the infrastructure). This structure includes (a) the "material forces of production," that is, the labour and means of production, and (b) the overall "relations of production," or the social and political arrangements that regulate production and distribution. Although Marx stated that there is a correspondence between the "material forces" of production and the indispensable "relations" of production, he never made himself clear on the nature of the correspondence, a fact that was to be the source of differing interpretations among his later followers.

2. Above the economic structure rises the superstructure consisting of legal and political "forms of social consciousness" that correspond to the economic structure. Marx says nothing about the nature of this correspondence between ideological forms and economic structure, except that through the ideological forms men become conscious of the conflict within the economic structure between the material forces of production and the existing relations of production expressed in the legal property relations. In other words, "The sum total of the forces of production accessible to men determines the condition of society" and is at the base of society. "The social structure and the state issue continually from the life processes of definite individuals . . . as they are *in reality*, that is acting and materially producing." The political relations that men establish among themselves are dependent on material production, as are the legal relations. This foundation of the social on the economic is not an incidental point: it colours Marx's whole analysis. It is found in *Das Kapital* as well as in *The German Ideology* and the *Economic and Philosophic Manuscripts of 1844*.

Analysis of society. To go directly to the heart of the work of Marx, one must focus on his concrete program for man. This is just as important for an understanding of Marx as are *The Communist Manifesto* and *Das Kapital*. Marx's interpretation of man begins with human need. "Man," he wrote in the *Economic and Philosophic Manuscripts of 1844*,

is first of all a *natural being*. As a natural being and a living natural being, he is endowed on the one hand with *natural powers, vital powers* . . . : these powers exist in him as aptitudes, instincts. On the other hand, as an objective, natural, physical, sensitive being, he is a *suffering*, dependent and limited being . . . that is, the *objects* of his instincts exist outside him, independent of him, but are the objects of his *need*, indispensable and essential for the realization and confirmation of his substantial powers.

The point of departure of human history is therefore living man, who seeks to satisfy certain primary needs. "The first historical fact is the production of the means to satisfy these needs." This satisfaction, in turn, opens the way for new needs. Human activity is thus essentially a struggle with nature that must furnish man with the means of satisfying his needs: drink, food, clothing, the development of his powers and then of his intellectual and artistic abilities. In this undertaking, man discovers himself as a productive being who humanizes himself by his labour. Furthermore, man humanizes nature while he naturalizes himself. By his creative activity, by his labour, he realizes his identity with the nature that he masters, while at the same time he achieves free consciousness. Born of nature man becomes fully human by opposing it. Becoming aware in his struggle against nature of what separates him from it, man finds the conditions of his fulfillment, of the realization of his true stature. The dawning of consciousness is inseparable from struggle. By appropriating all the creative energies, he discovers that "all that is called history is nothing else than the process of creating man through human labour, the becoming of nature for man. Man has thus evident and irrefutable proof of his own creation by himself." Understood in its universal dimension, human activity reveals that "for man, man is the supreme being." It is thus vain to speak of God, creation, and metaphysical problems. Fully naturalized, man is sufficient unto himself: he has recaptured the fullness of man in his full liberty.

Living in a capitalist society, however, man is not truly free. He is an alienated being; he is not at home in his world. The idea of alienation, which Marx takes from Hegel and Feuerbach, plays a fundamental role in the whole of his written work, starting with the writings of his youth and continuing through *Das Kapital*. In the *Economic and Philosophic Manuscripts* the alienation of labour is seen to spring from the fact that the more the worker produces the less he has to consume, and the more values he creates the more he devalues himself, because his product and his labour are estranged from him. The life of the worker depends on things that he has created but that are not his, so that, instead of finding his rightful existence through his labour, he loses it in this world of things that are external to him: no work, no pay. Under these conditions, labour denies the fullness of concrete man. "The generic being (*Gattungswesen*) of man, nature as well as his intellectual faculties, is transformed into a being which is alien to him, into a *means of his individual existence*." Nature, his body, his spiritual essence become alien to him. "Man is made alien to man." When carried to its highest stage of development, private property becomes "the product of alienated labour . . . the *means* by which labour alienates itself (and) the realization of this alienation." It is also at the same time "the tangible material expression of *alienated human life*."

Although there is no evidence that Marx ever disclaimed this anthropological analysis of alienated labour, starting with *The German Ideology*, the historical, social, and economic causes of the alienation of labour are given increasing emphasis, especially in *Das Kapital*. Alienated labour is seen as the consequence of market product, the division of labour, and the division of society into antagonistic classes. As producers in society, men create goods only by their labour. These goods are exchangeable. Their value is the average amount of social labour spent to produce them. The alienation of the worker takes on its full dimension in that system of market production in which part of the value of the goods produced by the worker is taken away from him and transformed into surplus value, which the capitalist privately appropriates. Market production also intensifies the alienation of labour by encouraging specialization, piecework, and the setting up of large enterprises. Thus the labour power of the worker is used along with that of others in a combination whose significance he is ignorant of, both individually and socially. In thus losing their quality as human products, the products of labour become fetishes, that is, alien and oppressive realities to which both the man who possesses them privately and the man who is deprived of them submit themselves. In the market economy, this submission to things is obscured by the fact that the exchange of goods is expressed in money.

This fundamental economic alienation is accompanied by secondary political and ideological alienations, which offer a distorted representation of and an illusory justification of a world in which the relations of men with one another are also distorted. The ideas that men form are closely bound up with their material activity and their material relations: "The act of making representations, of thinking, the spiritual intercourse of men, seem to be the direct emanation of their material relations." This is true of all human activity: political, intellectual, or spiritual. "Men produce their representations and their ideas, but it is as living men, men acting as they are determined by a definite development of their powers of production." Law, morality, metaphysics, and religion do not have a history of their own. "Men developing their material production modify together with their real existence their ways of thinking and the products of their ways of thinking." In other words, "It is not consciousness which determines existence, it is existence which determines consciousness."

In bourgeois, capitalist society man is divided into political citizen and economic man. This duality represents man's political alienation, which is further intensified by the functioning of the bourgeois state. From this study of society at the beginning of the 19th century, Marx came to see the state as the instrument through which the propertied class dominated other classes.

Ideological alienation, for Marx, takes different forms.

Man as an alienated being

Economic alienation

Man as a natural being

appearing in economic, philosophical, and legal theories. Marx undertook a lengthy critique of the first in *Das Kapital* and of the second in *The German Ideology*. But ideological alienation expresses itself supremely in religion. Taking up the ideas about religion that were current in left post-Hegelian circles, together with the thought of Feuerbach, Marx considered religion to be a product of man's consciousness. It is a reflection of the situation of a man who "either has not conquered himself or has already lost himself again" (man in the world of private property). It is "an opium for the people." Unlike Feuerbach, Marx believed that religion would disappear only with changes in society.

Analysis of the economy. Marx analyzed the market economy system in *Das Kapital*. In this work he borrows most of the categories of the classical English economists Smith and Ricardo but adapts them and introduces new concepts such as that of surplus value. One of the distinguishing marks of *Das Kapital* is that in it Marx studies the economy as a whole and not in one or another of its aspects. His analysis is based on the idea that man is a productive being and that all economic value comes from human labour. The system he analyzes is principally that of mid-19th-century England. It is a system of private enterprise and competition that arose in the 16th century from the development of sea routes, international trade, and colonialism. Its rise had been facilitated by changes in the forces of production (the division of labour and the concentration of workshops), the adoption of mechanization, and technical progress. The wealth of the societies that brought this economy into play had been acquired through an "enormous accumulation of commodities." Marx therefore begins with the study of this accumulation, analyzing the unequal exchanges that take place in the market.

According to Marx, if the capitalist advances funds to buy cotton yarn with which to produce fabrics and sells the product for a larger sum than he paid, he is able to invest the difference in additional production. "Not only is the value advance kept in circulation, but it changes in its magnitude, adds a plus to itself, makes itself worth more, and it is this movement that transforms it into capital." The transformation, to Marx, is possible only because the capitalist has appropriated the means of production, including the labour power of the worker. Now labour power produces more than it is worth. The value of labour power is determined by the amount of labour necessary for its reproduction or, in other words, by the amount needed for the worker to subsist and beget children. But in the hands of the capitalist the labour power employed in the course of a day produces more than the value of the sustenance required by the worker and his family. The difference between the two values is appropriated by the capitalist, and it corresponds exactly to the surplus value realized by capitalists in the market. Marx is not concerned with whether in capitalist society there are sources of surplus value other than the exploitation of human labour—a fact pointed out by Joseph Schumpeter (*Capitalism, Socialism, and Democracy*). He remains content with emphasizing this primary source:

Surplus value is produced by the employment of labour power. Capital buys the labour power and pays the wages for it. By means of his work the labourer creates new value which does not belong to him, but to the capitalist. He must work a certain time merely in order to reproduce the equivalent value of his wages. But when this equivalent value has been returned, he does not cease work, but continues to do so for some further hours. The new value which he produces during this extra time, and which exceeds in consequence the amount of his wage, constitutes surplus value.

Throughout his analysis, Marx argues that the development of capitalism is accompanied by increasing contradictions. For example, the introduction of machinery is profitable to the individual capitalist because it enables him to produce more goods at a lower cost, but new techniques are soon taken up by his competitors. The outlay for machinery grows faster than the outlay for wages. Since only labour can produce the surplus value from which profit is derived, this means that the capitalist's rate

of profit on his total outlay tends to decline. Along with the declining rate of profit goes an increase in unemployment. Thus, the equilibrium of the system is precarious, subject as it is to the internal pressures resulting from its own development. Crises shake it at regular intervals, preludes to the general crisis that will sweep it away. This instability is increased by the formation of a reserve army of workers, both factory workers and peasants, whose pauperization keeps increasing. "Capitalist production develops the technique and the combination of the process of social production only by exhausting at the same time the two sources from which all wealth springs: the earth and the worker." According to the Marxist dialectic, these fundamental contradictions can only be resolved by a change from capitalism to a new system.

Class struggle. Marx inherited the ideas of class and class struggle from Utopian socialism and the theories of Saint-Simon. These had been given substance by the writings of French historians such as Adolphe Thiers and François Guizot on the French Revolution of 1789. But unlike the French historians, Marx made class struggle the central fact of social evolution. "The history of all hitherto existing human society is the history of class struggles."

In Marx's view, the dialectical nature of history is expressed in class struggle. With the development of capitalism, the class struggle takes an acute form. Two basic classes, around which other less important classes are grouped, oppose each other in the capitalist system: the owners of the means of production, or bourgeoisie, and the workers, or proletariat. "The bourgeoisie produces its own grave-diggers. The fall of the bourgeoisie and the victory of the proletariat are equally inevitable" (*The Communist Manifesto*) because

the bourgeois relations of production are the last contradictory form of the process of social production, contradictory not in the sense of an individual contradiction, but of a contradiction that is born of the conditions of social existence of individuals; however, the forces of production which develop in the midst of bourgeois society create at the same time the material conditions for resolving this contradiction. With this social development the prehistory of human society ends.

When man has become aware of his loss, of his alienation, as a universal nonhuman situation, it will be possible for him to proceed to a radical transformation of his situation by a revolution. This revolution will be the prelude to the establishment of communism and the reign of liberty reconquered. "In the place of the old bourgeois society with its classes and its class antagonisms, there will be an association in which the free development of each is the condition for the free development of all."

But for Marx there are two views of revolution. One is that of a final conflagration, "a violent suppression of the old conditions of production," which occurs when the opposition between bourgeoisie and proletariat has been carried to its extreme point. This conception is set forth in a manner inspired by the Hegelian dialectic of the master and the slave, in *The Holy Family*. The other conception is that of a permanent revolution involving a provisional coalition between the proletariat and the petty bourgeoisie rebelling against a capitalism that is only superficially united. Once a majority has been won to the coalition, an unofficial proletarian authority constitutes itself alongside the revolutionary bourgeois authority. Its mission is the political and revolutionary education of the proletariat, gradually assuring the transfer of legal power from the revolutionary bourgeoisie to the revolutionary proletariat.

If one reads *The Communist Manifesto* carefully one discovers inconsistencies that indicate that Marx had not reconciled the concepts of catastrophic and of permanent revolution. Moreover, Marx never analyzed classes as specific groups of men opposing other groups of men. Depending on the writings and the periods, the number of classes varies; and unfortunately the pen fell from Marx's hand at the moment when, in *Das Kapital* (vol. 3), he was about to take up the question. Reading *Das Kapital*, one is furthermore left with an ambiguous impression with regard to the destruction of capitalism: will it be the result of the "general crisis" that Marx expects, or of the action of the conscious proletariat, or of both at once?

Theory
of surplus
value

The
crises of
capitalism

Two
views of
revolution

The contributions of Engels. Engels became a communist in 1842 and discovered the proletariat of England when he took over the management of the Manchester factory belonging to his father's cotton firm. In 1844, the year he began his close association and friendship with Marx, Engels was finishing his "Umrisse zu einer Kritik der Nationalökonomie" ("Outline of a Critique of Political Economy")—a critique of Smith, Ricardo, Mill, and J.B. Say. This remarkable study contained in seminal form the critique that Marx was to make of bourgeois political economy in *Das Kapital*. During the first years of his stay in Manchester, Engels observed carefully the life of the workers of that great industrial centre and described it in *Die Lage der arbeitenden Klassen in England* (*The Condition of the Working Class in England*), published in 1845 in Leipzig. This work was an analysis of the evolution of industrial capitalism and its social consequences. He collaborated with Marx in the writing of *The Holy Family*, *The German Ideology*, and *The Communist Manifesto*. The correspondence between them is of fundamental importance for the student of *Das Kapital*, for it shows how Engels contributed by furnishing Marx with a great amount of technical and economic data and by criticizing the successive drafts. This collaboration lasted until Marx's death and was carried on posthumously with the publication of the manuscripts left by Marx, which Engels edited, forming volumes 2 and 3 of *Das Kapital*. He also wrote various articles on Marx's work.

In response to criticism of Marx's ideas by a socialist named Eugen Dühring, Engels published several articles that were collected under the title *Herr Eugen Dührings Umwälzung der Wissenschaft*, which appeared in 1878 (*Herr Eugen Dühring's Revolution in Science* [*Anti-Dühring*]), and an unfinished work, *Dialektik und Natur* (1927; *Dialectics of Nature*), which he had begun around 1875–76. The importance of these writings to the subsequent development of Marxism can be seen from Lenin's observation that Engels "developed, in a clear and often polemical style, the most general scientific questions and the different phenomena of the past and present according to the materialist understanding of history and the economic theory of Karl Marx." But Engels was driven to simplify problems with a view to being pedagogical; he tended to schematize and systematize things as if the fundamental questions were settled. The connections that he thus established between some of Marx's governing ideas and some of the scientific ideas of his age gave rise to the notion that there is a complete Marxist philosophy. The idea was to play a significant role in the transition of Marxism from a "critique of daily life" to an integrated doctrine in which philosophy, history, and the sciences are fused.

Anti-Dühring is of fundamental importance for it constitutes the link between Marx and certain forms of modern Marxism. It contains three parts: Philosophy, Political Economy, and Socialism. In the first, Engels attempts to establish that the natural sciences and even mathematics are dialectical, in the sense that observable reality is dialectical: the dialectical method of analysis and thought is imposed on men by the material forces with which they deal. It is thus rightly applied to the study of history and human society. "Motion, in effect, is the mode of existence of matter," Engels writes. In using materialistic dialectic to make a critique of Dühring's thesis, according to which political forces prevail over all the rest in the molding of history, Engels provides a good illustration of the materialistic idea of history, which puts the stress on the prime role of economic factors as driving forces in history. The other chapters of the section Political Economy form a very readable introduction to the principal economic ideas of Marx: value (simple and complex), labour, capital, and surplus value. The section Socialism starts by formulating anew the critique of the capitalist system as it was made in *Das Kapital*. At the end of the chapters devoted to production, distribution, the state, the family, and education, Engels outlines what the socialist society will be like, a society in which the notion of value has no longer anything to do with the distribution of the goods produced because all labour "becomes at once and directly

social labour," and the amount of social labour that every product contains no longer needs to be ascertained by "a detour." A production plan will coordinate the economy. The division of labour and the separation of town and country will disappear with the "suppression of the capitalist character of modern industry." Thanks to the plan, industry will be located throughout the country in the collective interest, and thus the opposition between town and country will disappear—to the profit of both industry and agriculture. Finally, after the liberation of man from the condition of servitude in which the capitalist mode of production holds him, the state will also be abolished and religion will disappear by "natural death."

One of the most remarkable features of *Anti-Dühring* is the insistence with which Engels refuses to base socialism on absolute values. He admits only relative values, linked to historical, economic, and social conditions. Socialism cannot possibly be based on ethical principles: each epoch can only successfully carry out that of which it is capable. Marx had written this in his preface of 1859.

GERMAN MARXISM AFTER ENGELS

The work of Kautsky and Bernstein. The theoretical leadership after Engels was taken by Karl Kautsky, editor of the official organ of the German Social Democratic Party, *Die Neue Zeit*. He wrote *Karl Marx' ökonomische Lehren* (1887; *The Economic Doctrines of Karl Marx*), in which the work of Marx is presented as essentially an economic theory. Kautsky reduced the ideas of Marx and Marxist historical dialectic to a kind of evolutionism. He laid stress on the increasing pauperization of the working class and on the increasing degree of capitalist concentration. While opposing all compromise with the bourgeois state, he accepted the contention that the socialist movement should support laws benefiting the workers provided that they did not reinforce the power of the state. Rejecting the idea of an alliance between the working class and the peasantry, he believed that the overthrow of the capitalist state and the acquisition of political power by the working class could be realized in a peaceful way, without upsetting the existing structures. As an internationalist he supported peace, rejecting war and violence. For him, war was a product of capitalism. Such were the main features of "orthodox" German Marxism at the time when the "revisionist" theories of Eduard Bernstein appeared.

Bernstein created a great controversy with articles that he wrote in 1896 for *Die Neue Zeit*, arguing that Marxism needed to be revised. His divergence widened with the publication in 1899 of *Die Voraussetzungen des Sozialismus und die Aufgaben der Sozialdemokratie* (*Evolutionary Socialism*), to which rejoinders were made by Kautsky in *Bernstein und das Sozialdemokratische Programm: Eine Antikritik* (1899; "Bernstein and the Social Democratic Program") and the Polish-born Marxist Rosa Luxemburg in *Sozialreform oder Revolution* (*Reform or Revolution*), both in 1899. Bernstein focused first of all upon the labour theory of value. Along with the economists of his time he considered it outdated, both in the form expounded by British classical economists and as set forth in *Das Kapital*. He argued, moreover, that class struggle was becoming less rather than more intense, for concentration was not accelerating in industry as Marx had forecast, and in agriculture it was not increasing at all. Bernstein demonstrated this on the basis of German, Dutch, and English statistical data. He also argued that cartels and business syndicates were smoothing the evolution of capitalism, a fact that cast doubt on the validity of Marx's theory of capitalistic crises. Arguing that quite a few of Marx's theories were not scientifically based, Bernstein blamed the Hegelian and Ricardian structure of Marx's work for his failure to take sufficient account of observable reality.

To this, Kautsky replied that, with the development of capitalism, agriculture was becoming a sector more and more dependent on industry, and that in addition an industrialization of agriculture was taking place. Luxemburg took the position that the contradictions of capitalism did not cease to grow with the progress of finance capitalism and the exploitation of the colonies, and that these contradictions were leading to a war that would give the

Engels' critique of Dühring

Doctrinal disputes

proletariat its opportunity to assume power by revolutionary means.

The radicals. One of the most divisive questions was that of war and peace. This was brought to the fore at the outbreak of World War I, when Social Democratic deputies in the German Reichstag voted for the financing of the war. Among German Marxists who opposed the war were Karl Liebknecht and Luxemburg. Liebknecht was imprisoned in 1916 for agitating against the war. On his release in 1918 he took the leadership of the Spartacist movement, which was later to become the Communist Party of Germany. Luxemburg had also been arrested for her antimilitary activities. In addition to her articles, signed Junius, in which she debated with Lenin on the subject of World War I and the attitude of the Marxists toward it (published in 1916 as *Die Krise der Sozialdemokratie* [*The Crisis in the German Social-Democracy*]), she is known for her book *Die Akkumulation des Kapitals* (1913; *The Accumulation of Capital*). In this work she returned to Marx's economic analysis of capitalism, in particular the accumulation of capital as expounded in volume 2 of *Das Kapital*. There she found a contradiction that had until then been unnoticed: Marx's scheme seems to imply that the development of capitalism can be indefinite, though elsewhere he sees the contradictions of the system as bringing about increasingly violent economic crises that will inevitably sweep capitalism away. Luxemburg concluded that Marx's scheme is oversimplified and assumes a universe made up entirely of capitalists and workers. If increases in productivity are taken into account, she asserted, balance between the two sectors becomes impossible; in order to keep expanding, capitalists must find new markets in noncapitalist spheres, either among peasants and artisans or in colonies and underdeveloped countries. Capitalism will collapse only when exploitation of the world outside it (the peasantry, colonies, etc.) has reached a limit. This conclusion has been the subject of passionate controversies.

The Austrians. The Austrian school came into being when Austrian socialists started publishing their works independently of the Germans; it can be dated from either 1904 (beginning of the *Marx-Studien* collection) or 1907 (publication of the magazine *Der Kampf*). The most important members of the school were Max Adler, Karl Renner, Rudolf Hilferding, Gustav Eckstein, Friedrich Adler, and Otto Bauer. The most eminent was Bauer, a brilliant theoretician whose *Die Nationalitätenfrage und die Sozialdemokratie* (1906; "The Nationalities Question and the Social Democracy") was critically reviewed by Lenin. In this work he dealt with the problem of nationalities in the light of the experience of the Austro-Hungarian Empire. He favoured the self-determination of peoples and emphasized the cultural elements in the concept of nationhood. Hilferding was finance minister of the German Republic after World War I in the Cabinets of the Social Democrats Gustav Stresemann (1923) and Hermann Müller (1928). He is known especially for his work *Das Finanzkapital* (1910), in which he maintained that capitalism had come under the control of banks and industrial monopolies. The growth of national competition and tariff barriers, he believed, had led to economic warfare abroad. Hilferding's ideas strongly influenced Lenin, who analyzed them in *Imperialism, the Highest Stage of Capitalism* (1916).

RUSSIAN AND SOVIET MARXISM

Das Kapital was translated into Russian in 1872. Marx kept up more or less steady relations with the Russian socialists and took an interest in the economic and social conditions of the tsarist empire. The man who originally introduced Marxism into Russia was Georgi Plekhanov, but the man who adapted Marxism to Russian conditions was Lenin.

Lenin. Vladimir Ilich Ulyanov, or Lenin, was born in 1870 at Simbirsk. He entered the University of Kazan to study law but was expelled the same year for participating in student agitation. In 1893 he settled in St. Petersburg and became actively involved with the revolutionary workers. With his pamphlet *What Is To Be Done?* (1902),

he specified the theoretical principles and organization of a Marxist party as he thought it should be constituted. He took part in the second Congress of the Russian Social-Democratic Workers' Party, which was held in Brussels and London (1903), and induced the majority of the Congress members to adopt his views. Two factions formed at the Congress: the Bolshevik (from the Russian word for "larger") with Lenin as the leader and the Menshevik (from the Russian word for "smaller") with Julius Martov at the head. The former wanted a restricted party of militants and advocated the dictatorship of the proletariat. The latter wanted a wide-open proletarian party, collaboration with the liberals, and a democratic constitution for Russia. In his pamphlet *One Step Forward, Two Steps Back* (1904), Lenin compared the organizational principles of the Bolsheviks to those of the Mensheviks. After the failure of the 1905 Russian revolution, he drew positive lessons for the future in *Two Tactics of Social Democracy in the Democratic Revolution*. He fiercely attacked the influence of Kantian philosophy on German and Russian Marxism in *Materialism and Empirio-criticism* (1908). In 1912 at the Prague Conference the Bolsheviks constituted themselves as an independent party. During World War I Lenin resided in Switzerland, where he studied Hegel's *Science of Logic* and the development of capitalism and carried on debates with Marxists like Rosa Luxemburg on the meaning of the war and the right of nations to self-determination. In 1915 at Zimmerwald, and in 1916 at Kiental, he organized two international socialist conferences to fight against the war. Immediately after the February 1917 revolution he returned to Russia, and in October the Bolshevik coup brought him to power.

The situation of Russia and the Russian revolutionary movement at the end of the 19th century and the beginning of the 20th led Lenin to diverge, in the course of his development and his analyses, from the positions both of "orthodox Marxism" and of "revisionism." He rediscovered the original thought of Marx by a careful study of his works, in particular *Das Kapital* and *The Holy Family*. He saw Marxism as a practical affair and tried to go beyond the accepted formulas to plan political action that would come to grips with the surrounding world.

As early as 1894, in his populist study *The Friends of the People*, Lenin took up Marx's distinction between the "material social relations" of men and their "ideological social relations." In Lenin's eyes the importance of *Das Kapital* was that "while explaining the structure and the development of the social formation seen *exclusively* in terms of its relations of production, (Marx) has nevertheless everywhere and always analyzed the superstructure which corresponds to these relations of production." In *The Development of Russian Capitalism* (1897-99) Lenin sought to apply Marx's analysis by showing the growing role of capital, in particular commercial capital, in the exploitation of the workers in the factories and the large-scale expropriation of the peasants. It was thus possible to apply to Russia the models developed by Marx for western Europe. At the same time Lenin did not lose sight of the importance of the peasant in Russian society. Although a disciple of Marx, he did not believe that he had only to repeat Marx's conclusions. He wrote:

We do not consider the theory of Marx to be a complete, immutable whole. We think on the contrary that this theory has only laid the cornerstone of the science, a science which socialists must further develop in all directions if they do not want to let themselves be overtaken by life. We think that, for the Russian socialists, an independent elaboration of the theory is particularly necessary.

Lenin laid great stress upon the dialectical method. In his early writings he defined the dialectic as "nothing more nor less than the method of sociology, which sees society as a living organism, in perpetual development (and not as something mechanically assembled and thus allowing all sorts of arbitrary combinations of the various social elements) . . ." (*The Friends of the People*, 1894). After having studied Hegel toward the end of 1914, he took a more activist view. Dialectic is not only evolution; it is praxis, leading from activity to reflection and from reflection to action.

Karl
Liebknecht
and Rosa
Luxem-
burg

Lenin's
adapta-
tion of
Marxism

A plan of
political
action

The dictatorship of the proletariat. Lenin also put much emphasis on the leading role of the party. As early as 1902 he was concerned with the need for a cohesive party with a correct doctrine, adapted to the exigencies of the period, which would be a motive force among the masses, helping to bring them to an awareness of their real situation. In *What Is To Be Done?* he called for a party of professional revolutionaries, disciplined and directed, capable of defeating the police; its aim should be to establish the dictatorship of the proletariat. In order to do this, he wrote in *Two Tactics of Social-Democracy in the Democratic Revolution*, it was necessary "to subject the insurrection of the proletarian and non-proletarian masses to our influence, to our direction, to use it in our best interests." But this was not possible without a doctrine: "Without revolutionary theory, no revolutionary movement." On the eve of the revolution of October 1917, in *The State and Revolution* he set forth the conditions for the dictatorship of the proletariat and the suppression of the capitalist state.

Lenin assigned major importance to the peasantry in formulating his program. It would be a serious error, he held, for the Russian revolutionary workers' movement to neglect the peasants. Even though it was clear that the industrial proletariat constituted the vanguard of the revolution, the discontent of the peasantry could be oriented in a direction favourable to the revolution by placing among the goals of the party the seizure of privately owned land. As early as 1903, at the third congress of the party, he secured a resolution to this effect. Thereafter, the dictatorship of the proletariat became the dictatorship of the proletariat and the peasantry. In 1917 he encouraged the peasants to seize land long before the approval of agrarian reform by the Constituent Assembly.

Among Lenin's legacies to Soviet Marxism was one that proved to be injurious to the party. This was the decision taken at his behest by the 10th congress of the party in the spring of 1921, while the sailors were rebelling at Kronstadt and the peasants were growing restless in the countryside, to forbid all factions, all factional activity, and all opposition political platforms within the party. This decision had grave consequences in later years when Stalin used it against his opponents.

Stalin. It is Joseph Stalin who codified the body of ideas that, under the name of Marxism-Leninism, has constituted the official doctrine of the Soviet and eastern European communist parties. Stalin was a man of action in a slightly different sense than was Lenin. Gradually taking over power after Lenin's death in 1924, he pursued the development of the Soviet Union with great vigour. By practicing Marxism, he assimilated it, at the same time simplifying it. Stalin's Marxism-Leninism rests on the dialectic of Hegel, as set forth in *A Short History of the Communist Party of the Soviet Union* (1938), and on a materialism that can be considered roughly identical to that of Feuerbach. His work *Problems of Leninism*, which appeared in 11 editions during his lifetime, sets forth an ideology of power and activism that rides roughshod over the more nuanced approach of Lenin.

Soviet dialectical materialism can be reduced to four laws: (1) History is a dialectical development. It proceeds by successive phases that supersede one another. These phases are not separate, any more than birth, growth, and death are separate. Though it is true that phase B necessarily negates phase A, it remains that phase B was already contained in phase A and was initiated by it. The dialectic does not regard nature as an accidental accumulation of objects, of isolated and independent phenomena, but as a unified, coherent whole. Furthermore, nature is perpetually in movement, in a state of unceasing renewal and development, in which there is always something being born and developing and something disintegrating and disappearing. (2) Evolution takes place in leaps, not gradually. (3) Contradictions must be made manifest. All phenomena contain in themselves contradictory elements. "Dialectic starts from the point of view that objects and natural phenomena imply internal contradictions, because they all have a positive and a negative side." These contradictory elements are in perpetual struggle: it is this

struggle that is the "internal content of the process of development," according to Stalin. (4) The law of this development is economic. All other contradictions are rooted in the basic economic relationship. A given epoch is entirely determined by the relations of production existing among men. They are social relations; relations of collaboration or mutual aid, relations of domination or submission; and finally, transitory relations that characterize a period of passage from one system to another. "The history of the development of society is, above all, the history of the development of production, the history of the modes of production which succeed one another through the centuries."

From these principles may be drawn the following inferences, essential for penetrating the workings of Marxist-Leninist thought and its application. No natural phenomenon, no historical or social situation, no political fact, can be considered independently of the other facts or phenomena that surround it; it is set within a whole. Since movement is the essential fact, one must distinguish between what is beginning to decay and what is being born and developing. Since the process of development takes place by leaps, one passes suddenly from a succession of slow quantitative changes to a radical qualitative change. In the social or political realm, these sudden qualitative changes are revolutions, carried out by the oppressed classes. One must follow a frankly proletarian-class policy that exposes the contradictions of the capitalist system. A reformist policy makes no sense. Consequently (1) nothing can be judged from the point of view of "eternal justice" or any other preconceived notion and (2) no social system is immutable. To be effective, one must not base one's action on social strata that are no longer developing, even if they represent for the moment the dominant force, but on those that are developing.

Stalin's materialist and historical dialectic differs sharply from the perspective of Karl Marx. In *The Communist Manifesto* Marx applied the materialist dialectic to the social and political life of his time. In the chapter entitled "Bourgeois and Proletarians," he studied the process of the growth of the revolutionary bourgeoisie within feudal society, then the genesis and the growth of the proletariat within capitalism, placing the emphasis on the struggle between antagonistic classes. To be sure, he connected social evolution with the development of the forces of production. What counted for him, however, was not only the struggle but also the birth of consciousness among the proletariat. "As to the final victory of the propositions put forth in the *Manifesto*, Marx expected it to come primarily from the intellectual development of the working class, necessarily the result of common action and discussion" (Engels, preface to the republication of *The Communist Manifesto*, May 1, 1890).

The result of Stalin's dialectic, however, was what he called revolution from above, a dictatorial policy to increase industrialization and collectivize agriculture based upon ruthless repression and a strong centralization of power. For Stalin what counted was the immediate goal, the practical result. The move was from a dialectic that emphasized both the objective and the subjective to one purely objective, or more exactly, objectivist. Human actions are to be judged not by taking account of the intentions of the actor and their place in a given historical web but only in terms of what they signify objectively at the end of the period considered.

Trotskyism. Alongside Marxism-Leninism as expounded in the former Soviet Union, there arose another point of view expressed by Stalin's opponent Leon Trotsky and his followers. Trotsky played a leading role in both the Russian Revolution of 1905 and that of 1917. After Lenin's death he fell out with Stalin. Their conflict turned largely upon questions of policy, both domestic and foreign. In the realm of ideas, Trotsky held that a revolution in a backward, rural country could be carried out only by the proletariat. Once in power the proletariat must carry out agrarian reform and undertake the accelerated development of the economy. The revolution must be a socialist one, involving the abolition of the private ownership of the means of production, or else it will fail.

Differences between Marx's ideas and Stalin's

Stalin and dialectical materialism

But the revolution cannot be carried out in isolation, as Stalin maintained it could. The capitalist countries will try to destroy it; moreover, to succeed the revolution must be able to draw upon the industrial techniques of the developed countries. For these reasons the revolution must be worldwide and permanent, directed against the liberal and nationalist bourgeoisie of all countries and using local victories to advance the international struggle.

Tactically, Trotsky emphasized the necessity of finding or creating a revolutionary situation, of educating the working class in order to revolutionize it, of seeing that the party remained open to the various revolutionary tendencies and avoided becoming bureaucratized, and finally, when the time for insurrection comes, of organizing it according to a detailed plan.

VARIANTS OF MARXISM

Maoism. When the Chinese Communists took power in 1948, they brought with them a new kind of Marxism that came to be called Maoism after their leader Mao Zedong. The thought of Mao must always be seen against the changing revolutionary reality of China from 1930 onward. His thought was complex, a Marxist type of analysis combined with the permanent fundamentals of Chinese thought and culture.

One of its central elements has to do with the nature and role of contradictions in socialist society. For Mao, every society, including socialist (communist) society, contained "two different types of contradictions": (1) antagonistic contradictions—contradictions between us (the people) and our enemies (the Chinese bourgeoisie faithful), between the imperialist camp and the socialist camp, and so forth—which are resolved by revolution, and (2) nonantagonistic contradictions—between the government and the people under a socialist regime, between two groups within the Communist Party, between one section of the people and another under a communist regime, and so forth—which are resolved by vigorous fraternal criticism and self-criticism.

The notion of contradiction is specific to Mao's thought in that it differs from the conceptions of Marx or Lenin. For Mao, in effect, contradictions were at the same time universal and particular. In their universality, one must seek and discover what constitutes their particularity: every contradiction displays a particular character, depending on the nature of things and phenomena. Contradictions have alternating aspects—sometimes strongly marked, sometimes blurred. Some of these aspects are primary, others secondary. It is important to define them well, for if one fails to do so, the analysis of the social reality and the actions that follow from it will be mistaken. This is quite far from Stalinism and dogmatic Marxism-Leninism.

Another essential element of Mao's thought, which must be seen in the context of revolutionary China, is the notion of permanent revolution. It is an old idea advocated in different contexts by Marx, Lenin, and Trotsky but lacking, in Mao's formulation, the international dimension espoused by his predecessors. For Mao it followed from his ideas about the struggle of man against nature (held from 1938, at least); the campaigns for the rectification of thought (1942, 1951, 1952); and the necessity of struggling against bureaucracy, wastage, and corruption in a country of 600,000,000 to 700,000,000 inhabitants, where very old civilizations and cultures still permeated both the bourgeois classes and the peasantry, where bureaucracy was thoroughly entrenched, and where the previous society was extremely corrupt. It arose from Mao's conviction that the rhythm of the revolution must be accelerated. This conviction appeared in 1957 in his speeches and became manifest in 1958 in the "Great Leap Forward," followed in 1966 by the Cultural Revolution.

Mao's concept of permanent revolution rests upon the existence of nonantagonistic contradictions in the China of today and of tomorrow. Men must be mobilized into a permanent movement in order to carry forward the revolution and to prevent the ruling group from turning bourgeois (as he perceived it had in the Soviet Union). It is necessary to shape among the masses a new vision of the world by tearing them from their passivity and their

century-old habits. This is the background of the Cultural Revolution that began in 1966, following previous campaigns but differing from them in its magnitude and, it would seem, in the mobilization of youth against the cadres of the party. In these campaigns Mao drew upon his past as a revolutionary Marxist peasant leader, from his life in the red military and peasant bases and among the Red Guards of Yen-an, seeking in his past experience ways to mobilize the whole Chinese population against the dangers—internal and external—that confronted it in the present.

The distinguishing characteristic of Maoism is that it represents a peasant type of Marxism, with a principally rural and military outlook. While basing himself on Marxism-Leninism, adapted to Chinese requirements, Mao was rooted in the peasant life from which he himself came, in the revolts against the warlords and the bureaucrats that have filled the history of China. By integrating this experience into a universal vision of history, Mao gave it a significance that flows beyond the provincial limits of China.

In his effort to remain close to the Chinese peasant masses, Mao drew upon an idea of nature and a symbolism found in popular Chinese Taoism, though transformed by his Marxism. It can be seen in his many poems, which were written in the classical Chinese style. This idea of nature is accompanied in his written political works by the Promethean idea of man struggling in a war against nature, a conception in his thought that goes back at least to 1938 and became more important after 1955 as the rhythm of the revolution accelerated.

Marxism in Cuba. The Marxism of Fidel Castro expresses itself as a rejection of injustice in any form—political, economic, or social. In this sense it is related to the liberal democracy and Pan-Americanism of Simón Bolívar in Latin America during the 19th century. In its liberalism, Castro's early socialism resembled the various French socialisms of the first half of the 19th century. Only gradually did Castroism come to identify itself with Marxism-Leninism, although from the very beginning of the Cuban revolution Castro revealed his attachment to certain of Marx's ideas. Castro's Marxism rejects some of the tenets and practices of official Marxism-Leninism: it is outspoken against dogmatism, bureaucracy, and sectarianism. In one sense, Castroism is a Marxist-Leninist "heresy." It exalts the ethos of guerrilla revolution over party politics. At the same time it aims to apply a purer Marxism to the conditions of Cuba: alleged American imperialism, a single-crop economy, a low initial level of political and economic development. One may call it an attempt to realize a synthesis of Marxist ideas and the ideas of Bolívar.

In the ideological and political conflicts that divide the communist world, Castroism takes a more or less unengaged position. Castro is above all a nationalist and only after that a Marxist.

Marxism in the Third World. The development of Marxist variants in the Third World has been primarily influenced by the undeveloped industrial state and the former colonial status of the nations in question. In the traditional Marxist view the growth of capitalism is seen as a step necessary for the breakup of precapitalist peasant society and for the rise of the revolutionary proletariat class. Some theorists believe, however, that capitalism introduced by imperialist rather than indigenous powers sustains rather than destroys the feudal structure of peasant society and promotes underdevelopment because resources and surplus are usurped by the colonial powers. Furthermore, the revolutionary socialist movement becomes subordinate to that of national liberation, which violates Marx's theory of class struggle by uniting all indigenous classes in the common cause of anti-imperialism. For these reasons, many Third World countries have chosen to follow the Maoist model, with its emphasis on agrarian revolution against feudalism and imperialism, rather than the old Soviet one. Another alternative, one specific to the Third World, also exists. This policy bypasses capitalism and depends upon the established strength of other communist countries for support against imperialism.

Mao's theory of permanent revolution

Castroism

Marxism in the West. There are two main forms of Marxism in the West: that of the traditional communist parties and the more diffuse "New Left" form, which has come to be known as "Western Marxism." In general, the success of western European communist parties had been hindered by their perceived allegiance to the old Soviet authority rather than their own countries; the secretive, bureaucratic form of organization they inherited from Lenin; the ease with which they became integrated into capitalist society; and their consequent fear of compromising their principles by sharing power with bourgeois parties. The Western parties basically adhered to the policies of Soviet Marxism until the 1970s, when they began to advocate Eurocommunism, a moderate version of communism that they felt would broaden their base of appeal beyond the working class and thus improve their chances for political success. As described by Enrico Berlinguer, Georges Marchais, and Santiago Carrillo, the leaders of the Italian, French, and Spanish communist parties, respectively, Eurocommunism favoured a peaceful, democratic approach to achieving socialism, encouraged making alliances with other political parties, guaranteed civil liberties, and renounced the central authority of the Soviet party. By the 1980s Eurocommunism had largely been abandoned as unsuccessful, and communist parties in advanced capitalist nations returned to orthodox Marxism-Leninism despite the concomitant problems.

Western Marxism, however, can be seen as a repudiation of Marxism-Leninism, although, when it was first formulated in the 1920s, its proponents believed they were loyal to the dominant Soviet Communist Party. Prominent figures in the evolution of Western Marxism include the central Europeans György Lukács, Karl Korsch, and Lucien Goldmann; Antonio Gramsci of Italy; the German theorists who constituted the Frankfurt school, especially Max Horkheimer, Theodor Adorno, Herbert Marcuse, and Jürgen Habermas; and Henri Lefebvre, Jean-Paul Sartre, and Maurice Merleau-Ponty of France.

Western Marxism has been shaped primarily by the failure of the socialist revolution in the Western world. Western Marxists were concerned less with the actual political or economic practice of Marxism than with its philosophical interpretation, especially in relation to cultural and historical studies. In order to explain the inarguable success of capitalist society, they felt they needed to explore and understand non-Marxist approaches and all aspects of bourgeois culture. Eventually, they came to believe that traditional Marxism was not relevant to the reality of modern Western society.

Marx had predicted that revolution would succeed in Europe first, but, in fact, the Third World has proved more responsive. Orthodox Marxism also championed the technological achievements associated with capitalism, viewing them as essential to the progress of socialism. Experience showed the Western Marxists, however, that technology did not necessarily produce the crises Marx described and did not lead inevitably to revolution. In particular they disagreed with the idea, originally emphasized by Engels, that Marxism is an integrated, scientific doctrine that can be applied universally to nature; they viewed it as a critique of human life, not an objective, general science. Disillusioned by the terrorism of the Stalin era and the bureaucracy of the Communist Party system, they advocated the idea of government by workers' councils, which they believed would eliminate professional politicians and would more truly represent the interests of the working class. Later, when the working class appeared to them to be too well integrated into the capitalist system, the Western Marxists supported more anarchistic tactics. In general, their views are more in accord with those found in Marx's early, humanist writings rather than with his later, dogmatic interpretations.

Western Marxism has found support primarily among intellectuals rather than the working class, and orthodox Marxists have judged it impractical. Nevertheless, the Western Marxists' emphasis on Marx's social theory and their critical assessment of Marxist methodology and ideas have coloured the way even non-Marxists view the world.

(H.C./D.T.McL./Ed.)

MAJOR WORKS

Misère de la philosophie (1847; *The Poverty of Philosophy*, 1900); *Manifest der Kommunistischen Partei* (1848; *Manifesto of the Communists*, 1883); *Die Klassenkämpfe in Frankreich 1848 bis 1850* (1850; *The Class Struggles in France, 1848 to 1850*, 1924); *Der 18te Brumaire des Louis Napoleon* (1852; *The Eighteenth Brumaire of Louis Bonaparte*, 1898); *Zur Kritik der politischen Ökonomie* (1859; *A Contribution to the Critique of Political Economy*, 1904); and *Das Kapital* (vol. 1, 1867; vol. 2-3 published by Engels in 1885 and 1894; *Capital: A Critical Analysis of Capitalist Production*, vol. 1 trans. by Samuel Moore and Edward Aveling, 1886; vol. 2-3 trans. by Ernest Untermann, 1907 and 1909).

Recommended later translations of these works include *Manifesto of the Communist Party*, trans. by Samuel Moore (1888, reprinted 1952); *The Communist Manifesto of Karl Marx and Friedrich Engels*, with an introduction and explanatory notes by D. Ryazanoff, trans. by Eden Paul and Cedar Paul (1930); and *The Eighteenth Brumaire of Louis Bonaparte*, trans. by Eden Paul and Cedar Paul (1926). Selections from Marx's writings are available in the following: David McLellan (ed.), *Selected Writings* (1977); Robert C. Tucker (ed.), *The Marx-Engels Reader*, 2nd ed. (1978); and Jon Elster (ed.), *Karl Marx: A Reader* (1986).

A major English-language edition, *Karl Marx, Frederick Engels: Collected Works*, trans. by Richard Dixon et al. (1975-), is in progress. Planned to consist of 50 volumes and to include the correspondence, it is being prepared by an international editorial committee. Forty-one volumes had been published by 1992.

BIBLIOGRAPHY

Marx. The most comprehensive biography of Marx is DAVID MCLELLAN, *Karl Marx: His Life and Thought* (1973, reissued 1987). The classic biography of Marx, somewhat too partisan in his favour, is FRANZ MEHRING, *Karl Marx: The Story of His Life* (1935, reissued 1981; originally published in German, 1918). Marx's personal life is discussed in SAUL K. PADOVER, *Karl Marx, an Intimate Biography* (1978). JERROLD SEIGEL, *Marx's Fate: The Shape of a Life* (1978), is a psychological biography. YVONNE KAPP, *Eleanor Marx*, 2 vol. (1972-76), contains informative material on Marx's family life. Two good short biographies are ISAJAH BERLIN, *Karl Marx: His Life and Environment*, 4th ed. (1978, reprinted with corrections, 1982); and WERNER BLUMENBERG, *Portrait of Marx: An Illustrated Biography* (1972; originally published in German, 1962).

For introductory analysis and commentaries of Marx's works, see DAVID MCLELLAN, *The Thought of Karl Marx*, 2nd ed. (1980); BRUCE MAZLISH, *The Meaning of Karl Marx* (1984); W.A. SUCHTING, *Marx, an Introduction* (1983), and *Marx and Philosophy: Three Studies* (1986); and RICHARD SCHMITT, *Introduction to Marx and Engels: A Critical Reconstruction* (1987). TERRELL CARVER, *A Marx Dictionary* (1987), provides brief definitions of Marxian concepts without interpretative deviations from the original. For more detailed studies, see ROMAN ROSDOLSKY, *The Making of Marx's 'Capital'* (1977, reissued 1980; originally published in German, 1968); DEREK SAYER, *Marx's Method: Ideology, Science and Critique in Capital*, 2nd ed. (1983); ROBERT PAUL WOLFF, *Understanding Marx: A Reconstruction and Critique of Capital* (1984); HAL DRAPER, *Karl Marx's Theory of Revolution*, 3 vol. (1977-86); D. ROSS GANDY, *Marx and History: From Primitive Society to the Communist Future* (1979); MURRAY WOLFSON, *Marx: Economist, Philosopher, Jew: Steps in the Development of a Doctrine* (1982); DANIEL LITTLE, *The Scientific Marx* (1986); THOMAS SOWELL, *Marxism: Philosophy and Economics* (1985); and JOHN CUNNINGHAM WOOD (ed.), *Karl Marx's Economics: Critical Assessments*, 4 vol. (1987). A particularly acute summary of the difficulties in Marx's work is JON ELSTER, *Making Sense of Marx* (1985). Many monographs explore Marx's political and ideological development: CAROL C. GOULD, *Marx's Social Ontology: Individuality and Community in Marx's Theory of Social Reality* (1978, reprinted 1980), a study based on Marx's *Grundrisse*; RICHARD E. OLSEN, *Karl Marx* (1978); S.S. PRAWER, *Karl Marx and World Literature* (1976); PAUL THOMAS, *Karl Marx and the Anarchists* (1980); and ALLEN W. WOOD, *Karl Marx* (1981).

Marxism. Good introductions to the study of Marxism include LESZEK KOLAKOWSKI, *Main Currents of Marxism: Its Rise, Growth, and Dissolution*, 3 vol. (1978, reprinted 1981; originally published in Polish, 1976-78); GEORGE LICHTHEIM, *Marxism: An Historical and Critical Study*, 2nd ed. (1964, reprinted 1982); and DAVID MCLELLAN, *Marxism After Marx* (1979, reissued 1981), which contains an extensive bibliography. Some important analyses are assembled in DAVID MCLELLAN (ed.), *Marxism: Essential Writings* (1988). Studies of Marxism as a sociological doctrine may be found in KARL KORSCH, *Karl Marx* (1938, reissued 1963); HENRI LEBEVRE, *The Sociology of Marx* (1968, reprinted 1982; originally published in French, 1966);

and SIDNEY HOOK, *Towards the Understanding of Karl Marx* (1933). Developments in Marxism as a political theory are discussed in ALFRED SCHMIDT, *History and Structure: An Essay on Hegelian-Marxist and Structuralist Theories of History* (1981; originally published in German, 1971); DAVID RUBINSTEIN, *Marx and Wittgenstein: Social Praxis and Social Explanation* (1981); TOM ROCKMORE, *Fichte, Marx, and the German Philosophical Tradition* (1980); S.H. RIGBY, *Marxism and History: A Critical Introduction* (1987); and PAUL PHILLIPS, *Marx and Engels on Law and Laws* (1980). Specialized studies include STANLEY MOORE, *Marx on the Choice Between Socialism and Communism* (1980); JOSÉ PORFIRIO MIRANDA, *Marx Against the Marxists: The Christian Humanism of Karl Marx* (1980; originally published in Spanish, 1978); RALPH MILIBAND, *Marxism and Politics* (1977), including a discussion of the applicability of Marxism to contemporary politics in the Third World and communist countries; ROBERT L. HEILBRONER, *Marxism, For and Against* (1980); ANTHONY GIDDENS, *A Contemporary Critique of Historical Materialism*, 2 vol. (1981–85), an alternative, based on anthropological research, to the Marxist idea that all history has been the history of class struggle; MAURICE GODELIER, *Perspectives in Marxist Anthropology* (1977; originally published in French, 1973), presenting the contrasting view that classical Marxism may provide a methodology for analysis of empirical data in history and anthropology; and IAN CUMMINS, *Marx, Engels, and National Movements* (1980).

A. JAMES GREGOR, *A Survey of Marxism: Problems in Philosophy and the Theory of History* (1965), emphasizes philosophical problems in lieu of political or economic ones. The outstanding work on Marxist ethics is EUGENE KAMENKA, *The Ethical Foundations of Marxism*, 2nd ed. (1972). See also HUGO MEYNELL, *Freud, Marx, and Morals* (1981); and GARY NELSON and LAWRENCE GROSSBERG (eds.), *Marxism and Interpretation of Culture* (1988).

DAVID HOROWITZ (ed.), *Marx and Modern Economics* (1968), is an excellent collection of essays by leading economic theorists. Other treatments of Marxist economics worth consulting are PAUL M. SWEETZ, *The Theory of Capitalist Development: Principles of Marxian Political Economy* (1942, reissued 1970); and JOHN STRACHEY, *The Nature of Capitalist Crisis* (1935). The place of Marxist thought in the intellectual history of the 20th century is assessed in JACK LINDSAY, *The Crisis in Marxism* (1981); ANTHONY BREWER, *Marxist Theories of Imperialism: A Critical Survey* (1980); PERRY ANDERSON, *Considerations on Western Marxism* (1976); and WALTER L. ADAMSON, *Marx and the Disillusionment of Marxism* (1985).

An account of the historical development of Marxism can be found in HENRI CHAMBRE, *From Karl Marx to Mao Tse-Tung:*

A Systematic Survey of Marxism-Leninism (1963; originally published in French, 1959). GEORGE D.H. COLE, *A History of Socialist Thought*, 5 vol. in 7 (1953–65), presents a detailed study of the Marxist movement rather than the ideas; see especially vol. 2, *Socialist Thought: Marxism and Anarchism, 1850–1890*. TOM BOTTOMORE (ed.), *Interpretations of Marx* (1988), is an authoritative collection of essays.

The development and influence of Russian, Soviet, and eastern European Marxist theories is covered in a number of works by both Marxist and non-Marxist authors: HERBERT MARCUSE, *Soviet Marxism: A Critical Analysis* (1958, reprinted 1985); BERTRAM D. WOLFE, *Revolution and Reality: Essays on the Origin and Fate of the Soviet System* (1981); BARUCH KNEI-PAZ, *The Social and Political Thought of Leon Trotsky* (1978); UMBERTO MELOTTI, *Marx and the Third World* (1977, reprinted 1982; originally published in Italian, 1971); ADAM WESTOBY, *Communism Since World War II* (1981); and ERNEST MANDEL, *Revolutionary Marxism Today* (1979). Specialized studies include DONALD C. HODGES, *The Bureaucratization of Socialism* (1981); ROBERT J. ALEXANDER, *The Right Opposition: The Lovestoneites and the International Communist Opposition of the 1930's* (1981); ESTHER KINGSTON-MANN, *Lenin and the Problem of Marxist Peasant Revolution* (1983); BOGDAN SZAJKOWSKI (ed.), *Marxist Governments: A World Survey*, 3 vol. (1981); V. KUBÁLKOVÁ and A.A. CRUICKSHANK, *Marxism-Leninism and Theory of International Relations* (1980); HORACE B. DAVIS, *Toward a Marxist Theory of Nationalism* (1978); ISAAC DEUTSCHER, *Marxism in Our Time* (1971); JOHN P. BURKE, LAWRENCE CROCKER, and LYMAN H. LEGTERS (eds.), *Marxism and the Good Society* (1981), on Russia and China; JOHN G. GURLEY, *Challengers to Capitalism: Marx, Lenin, Stalin, and Mao*, 3rd ed. (1988); and NICHOLAS ABERCROMBIE, STEPHEN HILL, and BRYAN S. TURNER, *The Dominant Ideology Thesis* (1980), a critique of current Marxist thought. Two important critical studies are DAVID LANE, *The Socialist Industrial State: Towards a Political Sociology of State Socialism* (1976); and DONALD WILHELM, *Creative Alternatives to Communism: Guidelines for Tomorrow's World* (1977, reprinted 1981).

Of special reference interest are JOHN LACHS, *Marxist Philosophy: A Bibliographical Guide* (1967); HARRY G. SHAFFER, *Periodicals on the Socialist Countries and on Marxism: A New Annotated Index of English-Language Publications* (1977); J. WILCZYNSKI, *An Encyclopedic Dictionary of Marxism, Socialism and Communism* (1981); and ROBERT A. GORMAN (ed.), *Biographical Dictionary of Marxism* (1986), and *Biographical Dictionary of Neo-Marxism* (1985), a compendium providing information on practitioners of Marxism in more than 50 countries. (D.T.McL.)

Masks

Simply defined, a mask is a form of disguise. It is an object that is frequently worn over or in front of the face to hide the identity of a person and by its own features to establish another being. This essential characteristic of hiding and revealing personalities or moods is common to all masks. As cultural objects they have been used throughout the world in all periods since the Stone Age and have been as varied in appearance as in their use and symbolism.

This article deals with the general characteristics, functions, and forms of masks. It is divided into the following sections:

General characteristics	544
The making of masks	
The wearing of masks	
The role of the spectator	
Meaning and aesthetic response	
Preservation and collecting	
The functions and forms of masks	545
Social and religious uses	
Funerary and commemorative uses	
Therapeutic uses	
Festive uses	
Theatrical uses	
Bibliography	551

GENERAL CHARACTERISTICS

Masks have been designed in innumerable varieties, from the simplest of crude "false faces" held by a handle to complete head coverings with ingenious movable parts and hidden faces. Mask makers have shown great resourcefulness in selecting and combining available materials. Among the substances utilized are woods, metals, shells, fibres, ivory, clay, horn, stone, feathers, leather, furs, paper, cloth, and corn husks. Surface treatments have ranged

Carl Frank



Mask worn with costume: *makishi* dancer, a masked ancestral spirit who assists at initiation rites of the tribes of the northwestern region of Zambia.

from rugged simplicity to intricate carving and from polished woods and mosaics to gaudy adornments.

Masks generally are worn with a costume, often so complete that it entirely covers the body of the wearer. Fundamentally the costume completes the new identity represented by the mask, and usually tradition prescribes its appearance and construction to the same extent as the mask itself. Costumes, like the masks, are made of a great variety of materials, all of which have a symbolic connection with the mask's total imagery. Ideally the costume should be seen with the mask while the wearer is in action.

The morphological elements of the mask are with few exceptions derived from natural forms. Masks with human features are classified as anthropomorphic and those with animal characteristics as theriomorphic. In some instances, the mask form is a replication of natural features or closely follows the lineaments of reality, and in other instances it is an abstraction. Masks usually represent supernatural beings, ancestors, and fanciful or imagined figures and can also be portraits. The localization of a particular spirit in a specific mask must be considered a highly significant reason for its existence. The change in identity of the wearer for that of the mask is vital, for if the spirit represented does not reside in the image of the mask, the ritual petitions, supplications, and offerings made to it would be ineffectual and meaningless. The mask, therefore, most often functions as a means of contact with various spirit powers, thereby protecting against the unknown forces of the universe by prevailing upon their potential beneficence in all matters relative to life.

The making of masks. With few exceptions, masks have been made by professionals who were either expert in this particular craft or were noted sculptors or artisans. In societies in which masks of supernatural beings have played a significant ceremonial role, it is presumed that the spirit power of the created image usually is strongly felt by the artist. A primary belief involved in both the conception and the rendering of these objects was that spirit power dwelled in all organic and inorganic matter, and therefore the mask will contain the spirit power of whatever material was used to make it. This power is considered a volatile, active force that is surrounded by various taboos and restrictions for the protection of those handling it. Certain prescribed rituals frequently have to be followed in the process of the mask's creation. A spirit power is also often believed to inhabit the artist's tools so that even these have to be handled in a prescribed manner. As the form of the mask develops it is usually believed to acquire power increasingly in its own right, and again various procedures are prescribed to protect the craftsman and to ensure the potency of the object. If all the conventions have been adhered to, the completed mask, when worn or displayed, is regarded as an object suffused with great supernatural or spirit power. In some cultures it is believed that because of the close association between the mask maker and the spirit of the mask, the artist absorbs some of its magic power. A few West African tribal groups in Mali believe, in fact, that the creators of masks are even potentially capable of using the object's supernatural powers to cause harm to others.

Aesthetically, the mask maker has usually been restricted in the forms he can use since masks generally have a traditional imagery with formal conventions. If they are not followed, the artist can bring upon himself the severe censure of his social group and the displeasure or even wrath of the spirit power inherent in the mask. This requirement for accuracy, however, does not restrict artistic expressiveness. The mask maker can and does give his own creative interpretation to the traditionally prescribed general forms, attributes, and devices. The artist, in fact, is usually sought out as a maker of masks because of his

Anthropomorphic and theriomorphic masks

known ability to give a vitally expressive or an aesthetically pleasing presentation of the required image.

The wearing of masks. The wearer is also considered to be in direct association with the spirit force of the mask and is consequently exposed to like personal danger of being affected by it. For his protection, the wearer, like the mask maker, is required to follow certain sanctioned procedures in his use of the mask. In some respects he plays the role of an actor in cooperation or collaboration with the mask. Without his performing dance and posturing routines, which are often accompanied with certain sounds of music, the mask would remain a representation without a full life-force. The real drama and power of its form is the important contribution of the wearer. When he is attired in the mask, there is a loss of his previous identity and the assuming of a new one. Upon donning the mask, the wearer sometimes undergoes a psychic change and as in a trance assumes the spirit character depicted by the mask. Usually, however, the wearer skillfully becomes a "partner" of the character he is impersonating, giving to the mask not only an important spark of vitality by the light flashing from his own eyes but also bringing it alive by his movements and poses. But it would seem that the wearer often becomes psychologically completely attached to the character he is helping to create. He loses his own identity and becomes like an automaton, without his own will, which has become subservient to that of the personage of the mask. It appears, however, that at all times there remains some important, even if *sub rosa*, association between the mask and its wearer.

The role of the spectator. It is as consecrated objects imbued with supernatural power that masks are viewed by the spectators or participants at ceremonials where their presence is required. Whatever their specific identity may be, the masks usually refer back to early times, when their initial appearance occurred. This basic aspect of the mask is understood at least in essence by everyone. A paramount role of the mask is to give a sense of continuity between the present and the beginnings of time, a sense that is of vital importance for the integration of a culture with no written history. Psychologically the spectators become associated with the past through the spirit power of the mask, and this often leads the participants to a state of complete absorption or near-frenzy. This is not, however, a consistent reaction to masked ceremonials. That depends on the character whose presence the mask represents. In some cases, the spirit or supernatural being depicted is viewed with rejoicing and almost a familiarity, which leads to gaiety that has a cathartic aspect. Even so, the mask has a spirit content that is respected and revered, even if it is not showing a being with malignant potential. All of these forms have spirit and magical qualities and are thus esteemed as agents for the accomplishing of suprahuman acts.

Some masks, however, do represent malignant, evil, or potentially harmful spirits. These are often used to keep a required balance of power or a traditional social and political relationship of inherited positions within a culture. The characters depicted are also prescribed by tradition and enact roles to achieve the desired ends. The drama involving these masks is often associated with secret societies, especially in Africa, where the greatest range of mask forms and functions can be observed. These forms are often used in very restricted performances, where only select persons can view them. This is also true in other areas where masks are used, such as in Oceania, the Americas, and even in some of the folk mask rites still performed in Europe.

Meaning and aesthetic response. On the basis of present knowledge, it would appear that there is not or has not been any set response or reaction by any one of the three groups involved with the mask: the artist, the wearer, the spectator. There is, however, a reaction of a very particular kind common to every culture, a response such as awe, delight and pleasure, fear and even terror: these are as traditionally determined as the forms and costumes of the masks themselves. This is a learned and inherent pattern of conduct for each culture. Masks, therefore, that have a closely comparable appearance in several unrelated

groups in quite different parts of the world often have totally dissimilar meanings and functions. It is thus practically impossible to determine either the meaning or use of a mask by its appearance alone. For example, some masks in Africa, as well as in Oceania and East Asia, have such a grotesque or frightening appearance as to lead one to suspect that they represent evil spirits with an intent to terrorize the spectators; actually they may have the opposite character and function. The significance of masks can be determined only by reference to accounts or personal observations of the masks in the setting of their own culture.

The aesthetic effects of masks, on the other hand, since they derive from the forms and their disposition within the design, can readily be evaluated as art objects. But this evaluation is based on elements very different from those appraised within the mask's own culture. This is partly because the total artistic qualities of a mask derive both from its exterior forms and from its meaning and function within its cultural context. There exist, however, in all cultures criteria for determining the quality of objects as art. These criteria differ from one culture to another, and they may be known only from investigations carried out within the varying cultures.

Preservation and collecting. The preservation or disposal of masks is often decreed by tradition. Many masks and often their form and function are passed down through clans, families, special societies, or from individual to individual. They are usually spiritually reactivated or aesthetically restored by repainting and redecorating, without destroying the basic form and symbolism. In many instances, however, the mask is used only for one ceremony or occasion and then is discarded or destroyed, sometimes by burning.

The collecting of masks has largely been of recent origin. Not until the late 19th and early 20th century were they seriously appreciated as art objects or studied as cultural artifacts. Most masks have been obtained through archaeological excavations or in field expeditions, that is, in their place of origin.

THE FUNCTIONS AND FORMS OF MASKS

Masks are as extraordinarily varied in appearance as they are in function or fundamental meaning. Many masks are primarily associated with ceremonies that have religious and social significance or are concerned with funerary customs, fertility rites, or curing sickness. Other masks are used on festive occasions or to portray characters in a dramatic performance and in re-enactments of mythological events. Masks are also used for warfare and as protective devices in certain sports, as well as frequently being employed as architectural ornament.

Social and religious uses. Masks representing potentially harmful spirits were often used to keep a required balance of power or a traditional relationship of inherited positions within a culture. The forms of these masks invariably were prescribed by tradition, as were their uses. This type of mask was often associated with secret societies, especially in Africa, where the greatest range of types and functions can be found. They were also widely used among Oceanic peoples of the South Pacific and the American Indians and are even used in some of the folk rites still performed in Europe.

Masks have served an important role as a means of discipline and have been used to admonish women, children, and criminals. Common in China, Africa, Oceania, and North America, admonitory masks usually completely cover the features of the wearer. It is believed among some of the African Negro tribes that the first mask was an admonitory one. A child, repeatedly told not to, persisted in following its mother to fetch water. To frighten and discipline the child, the mother painted a hideous face on the bottom of her water gourd. Others say the mask was invented by a secret African society to escape recognition while punishing marauders. In New Britain, members of a secret terroristic society called the Dukduk appear in monstrous five-foot masks to police, to judge, and to execute offenders. Aggressive supernatural spirits of an almost demonic nature are represented by these masks, which are

Art
objects

Supernatural
and spirit
powers

Admonitory and
disciplinary masks



Admonitory and ancestor masks.

(Left) Female *tumbuan* mask of the Dukduk secret society, New Britain. The head is of sacking with tuft of feathers and skirt of leaves. Height 1.5 m. (Right) Ancestor mask from the middle Sepik River area of New Guinea, used in clan initiation rites. Upper portion is basketwork with faces modelled of red earth mixed with coconut oil, trimmed with shells, feathers, and a skirt of bark fibres. Height 1.79 m.

By courtesy of the Museum für Volkerkunde, Basel, Switzerland

constructed from a variety of materials, usually including tapa, or bark cloth, and the pith of certain reeds. These materials are painted in brilliant colours, with brick red and acid green predominating.

In many cultures throughout the world, a judge wears a mask to protect him from future recriminations. In this instance, the mask represents a traditionally sanctioned spirit from the past who assumes responsibility for the decision levied on the culprit.

Ancestor masks

Rituals, often nocturnal, by members of secret societies wearing ancestor masks are reminders of the ancient sanction of their conduct. In many cultures, these masked

ceremonials are intended to prevent miscreant acts and to maintain the circumscribed activities of the tribe. Along the Guinea coast of West Africa, for instance, many highly realistic masks represent ancestors who enjoyed specific cultural roles; the masks symbolize sanction and control when donned by the wearer. Among some of the Dan and Ngere tribes of Liberia and Ivory Coast, ancestor masks with generic features act as intermediaries for the transmission of petitions or offerings of respect to the gods. These traditional ancestral emissaries exert by their spirit power a social control for the community.

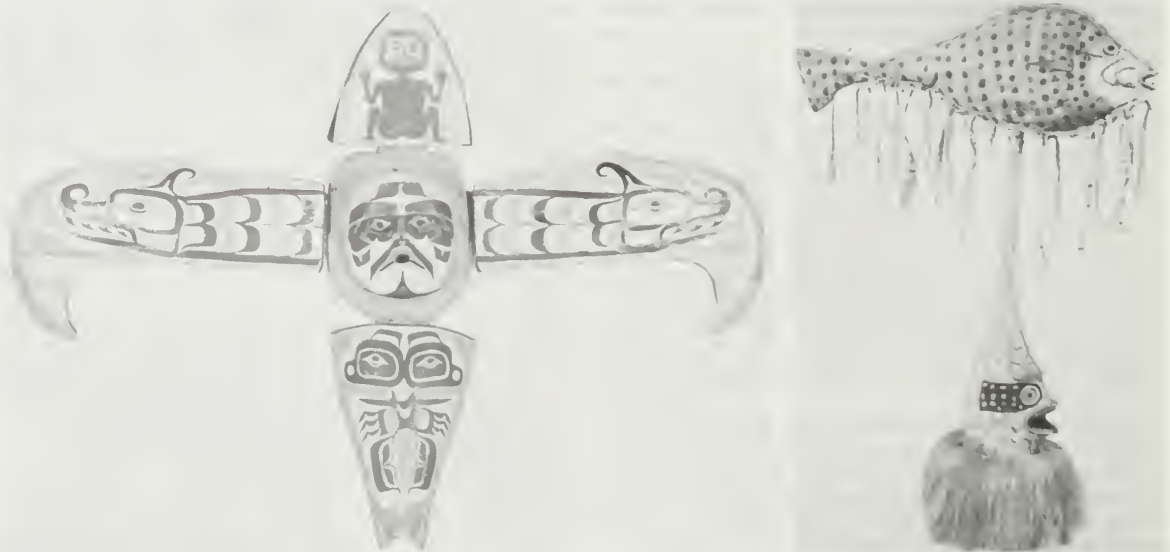
Particularly among Oceanic peoples, American Indians, and Negro tribes of Africa, certain times of the year are set aside to honour spirits or ancestors. Among nonliterate peoples who cannot record their own histories, masked rituals act as an important link between past and present, giving a sense of historic continuity that strengthens their social bond. On these occasions, masks usually recognizable as dead chieftains, relatives, friends, or even foes are worn or exhibited. Gifts are made to the spirits incarnated in the masks, while in other instances dancers wearing stylized mourning masks perform the prescribed ceremony.

In western Melanesia, the ancestral ceremonial mask occurs in a great variety of forms and materials. The Sepik River area in north central New Guinea is the source of an extremely rich array of these mask forms mostly carved in wood, ranging from small faces to large fantastic forms with a variety of appendages affixed to the wood, including shell, fibre, animal skins, seed, flowers, and feathers. These masks are highly polychromed with earth colours of red and yellow, lime white, and charcoal black. They often represent supernatural spirits as well as ancestors and therefore have both a religious and a social significance.

Members of secret societies usually conduct the rituals of initiation, when a young man is instructed in his future role as an adult and is acquainted with the rules controlling the social stability of the tribe. Totem and spiritualistic masks are donned by the elders at these ceremonies. Sometimes the masks used are reserved only for initiations. Among the most impressive of the initiation masks are the exquisitely carved human faces of west coast African Negro tribes. In western and central Zaire, in Africa, large, colourful helmetlike masks are used as a masquerading device when the youth emerges from the initiation area and is introduced to the villagers as an adult of the tribe. After a lengthy ordeal of teaching and initiation rites, for instance, a youth of the Pende tribe appears in a distinctive colourful mask indicative of his new role as an adult. The mask is later cast aside and replaced by a

Initiation masks

By courtesy of (left) the Brooklyn Museum, New York; (right) the Pitt Rivers Museum, Oxford, England



Totem masks

(Left) Thunderbird mask of the Kwakiutl Indians of the Northwest Coast of North America. The bird's beak and wings open to reveal a human face (mask shown open). Painted wood. Width (open) 1.82 m. (Right) Fish mask from the Orokol Bay area of New Guinea. Painted bark cloth over rattan frame with fringe of dried grass. Height 1.63 m.

Totem masks

small ivory duplicate, worn as a charm against misfortune and as a symbol of his manhood.

Believing everything in nature to possess a spirit, man found authority for himself and his family by identifying with a specific nonhuman spirit. He adopted an object of nature; then he mythologically traced his ancestry back to the chosen object; he preempted the animal as the emblem of himself and his clan. This is the practice of totem, which consolidates family pride and distinguishes social lines. Masks are made to house the totem spirit. The totem ancestor is believed actually to materialize in its mask; thus masks are of the utmost importance in securing protection and bringing comfort to the totem clan.

The Papuans of New Guinea build mammoth masks called *hevehe*, attaining 20 feet in height. They are constructed of a palm wood armature covered in bark cloth; geometric designs are stitched on with painted cane strips. These fantastic man-animal masks are given a frightening aspect. When they emerge from the men's secret clubhouse, they serve to protect the members of the clan. The so-called "totem" pole of the Alaskan and British Columbian Indian fulfills the same function. The African totem mask is often carved from ebony or other hard woods, designed with graceful lines and showing a highly polished surface. Animal masks, their features elongated and beautifully formalized, are common in western Africa. Dried grass, woven palm fibres, coconuts, and shells, as well as wood are employed in the masks of New Guinea, New Ireland, and New Caledonia. Represented are fanciful birds, fishes, and animals with distorted or exaggerated features.

The high priest and medicine man, or the shaman,

frequently had his own very powerful totem, in whose mask he could exorcise evil spirits, punish enemies, locate game or fish, predict the weather, and, most importantly, cure disease.

The Northwest Coast Indians of North America in particular devised mechanical masks with movable parts to reveal a second face—generally a human image. Believing that the human spirit could take animal form and vice versa, the makers of these masks fused man and bird or man and animal into one mask. Some of these articulating masks acted out entire legends as their parts moved.

Funerary and commemorative uses. In cultures in which burial customs are important, anthropomorphic masks have often been used in ceremonies associated with the dead and departing spirits. Funerary masks were frequently used to cover the face of the deceased. Generally their purpose was to represent the features of the deceased, both to honour them and to establish a relationship through the mask with the spirit world. Sometimes they were used to force the spirit of the newly dead to depart for the spirit world. Masks were also made to protect the deceased by frightening away malevolent spirits.

From the Middle Kingdom (c. 2040–1786 BC) to the 1st century AD, the ancient Egyptians placed stylized masks with generalized features on the faces of their dead. The funerary mask served to guide the spirit of the deceased back to its final resting place in the body. They were commonly made of cloth covered with stucco or plaster, which was then painted. For more important personages, silver and gold were used. Among the most splendid examples of the burial portrait mask is the one created c. 1350 BC for the pharaoh Tutankhamen. In Mycenaean tombs of

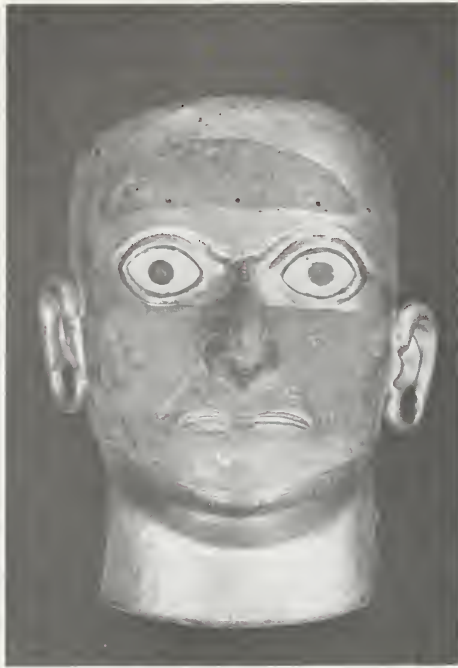
Burial masks

By courtesy of (top left) Museo Egizio, Turin, Italy, (top right) the National Archaeological Museum of Athens, TAP Service, (bottom left) the trustees of the British Museum, photographs, (top left) Foto Rosso, (bottom right) Nicholas Servian/Woodmansterne Limited



Funerary masks.

(Top left) Egyptian burial mask of the Ptolemaic period. Painted stucco over cloth. (Top right) Gold funeral mask placed over the face of an unknown Mycenaean ruler. Greece, 13th century BC. (Bottom left) Aztec skull mask from Mexico inlaid with turquoise and lignite with pyrite in the eye sockets. (Bottom right) Effigy mask of Edmund Sheffield, duke of Buckingham, 1735. Painted death mask with human hair, dressed in his peer's robes.



(Left) Metallic death mask with movable ears from the Moon Pyramid in the Moche Valley, Peru, 3rd–8th century AD. (Right) Ancient Mexican mask made of porphyry found north of Texcoco, Mexico, Teotihuacán civilization, 3rd–4th century AD.

By courtesy of (left) the Linden-Museum, Stuttgart, Germany (right) the American Museum of Natural History, New York City

c. 1400 BC, beaten gold portrait masks were found. Gold masks also were placed on the faces of the dead kings of Cambodia and Siam.

Use by the Incas

The mummies of Inca royalty wore golden masks. The mummies of lesser personages often had masks that were made of wood or clay. Some of these ancient Andean masks had movable parts, such as the metallic death mask with movable ears that was found in the Moon Pyramid at Moche, Peru. The ancient Mexicans made burial masks that seem to be generic representations rather than portraits of individuals.

In ancient Roman burials, a mask resembling the deceased was often placed over his face or was worn by an actor hired to accompany the funerary cortege to the burial site. In patrician families these masks or *imagines* were sometimes preserved as ancestor portraits and were displayed on ceremonial occasions. Such masks were usually modelled over the features of the dead and cast in wax. This technique was revived in the making of effigy masks for the royalty and nobility of Europe from

the late Middle Ages through the 18th century. Painted and with human hair, these masks were attached to a dummy dressed in state regalia and were used for display, processions, or commemorative ceremonials. From the 17th century to the 20th, death masks of famous persons became widespread among European peoples. With wax or liquid plaster of paris, a negative cast of the human face could be produced that in turn acted as a mold for the positive image, frequently cast in bronze. In the 19th century, life masks made in the same manner became popular. Another type of life mask had been produced in the Fayyum region of Egypt during the 1st and 2nd centuries AD. These were realistic portraits painted in encaustic on wood during a person's lifetime; when the person died, they were attached directly to the facial area on the mummy shroud.

Death and life masks

The skull mask is another form usually associated with funerary rites. The skull masks of the Aztecs, like their wooden masks, were inlaid with mosaics of turquoise and lignite, and the eye sockets were filled with pyrites.

By courtesy of (left) the trustees of the British Museum, (right) the American Museum of Natural History, New York City



Therapeutic masks

(Left) Disease devil mask (*rakasa*) from Sri Lanka worn to cure patients of deafness. (Right) False Face Society mask of the Iroquois Indians of North America.

Holes were customarily drilled in the back so the mask might be hung or possibly worn. In Melanesia, the skull of the deceased is often modelled over with clay, or resin and wax, and then elaborately painted with designs that had been used ceremonially by the deceased during his own lifetime.

Therapeutic uses. Masks have played an important part in magico-religious rites to prevent and to cure disease. In some cultures, the masked members of secret societies could drive disease demons from entire villages and tribes. Among the best known of these groups was the False Face Society of the North American Iroquois Indians. These professional healers performed violent pantomimes to exorcise the dreaded *Gahadogoka gogosa* (demons who plagued the Iroquois). They wore grimacing, twisted masks, often with long wigs of horsehair. Metallic inserts often were used around the eyes to catch the light of the campfire and the moon, emphasizing the grotesqueness of the mask.

Masks for protection from disease include the measles masks worn by Chinese children and the cholera masks worn during epidemics by the Chinese and Burmese. The disease mask is most developed among the Sinhalese in Sri Lanka (Ceylon), where 19 distinct *rākasa*, or disease devil masks, have been devised. These masks are of ferocious aspect, fanged, and with startling eyes. Gaudily coloured and sometimes having articulating jaws, they present a dragon-like appearance.

Masks have long been used in military connections. A war mask will have a malevolent expression or hideously fantastic features to instill fear in the enemy. The ancient Greeks and Romans used battle shields with grotesque masks or attached terrifying masks to their armour, as did the Chinese warrior. Grimacing *menpo*, or mask helmets, were used by Japanese samurai.

Many sports require the use of masks. Some of these are merely functional, protective devices such as the masks worn by fencers, baseball catchers, or even skiers. To protect their faces in sports events and tournaments of arms, horsemen of the Roman army attached highly decorative and symbolic masks to their helmets.

Perhaps the earliest use of masks was in connection with hunting. Disguise masks were seemingly used in the early Stone Age in stalking prey and later to house the slain animal's spirit in the hope of placating it. The traditional animal masks worn by the Altaic and Tungusic shamans in Siberia are strictly close to such prehistoric examples as the image of the so-called sorcerer in the Cave of Les Trois Frères in Ariège, France.

Since agricultural societies first appeared in prehistory, the mask has been widely used for fertility rituals. The Iroquois, for instance, used corn husk masks at harvest

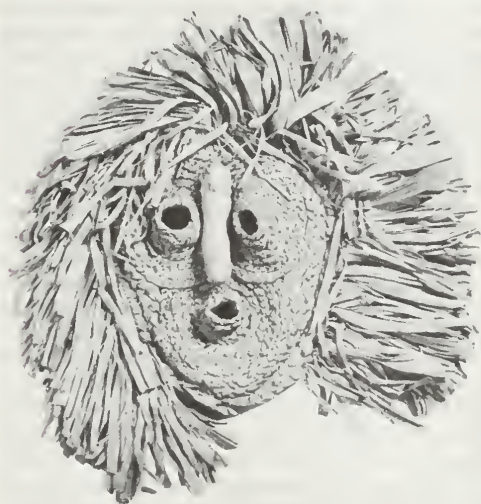
rituals to give thanks for and to achieve future abundance of crops. Perhaps the most renowned of the masked fertility rites held by American Indians are those still performed by the Hopi and Zuni Indians of the Southwest U.S. Together with masked dancers representing clouds, rain spirits, stars, Earth Mother, sky god, and others, the shaman takes part in elaborate ceremonies designed to assure crop fertility. Spirits called *kachinas*, who first brought rain to the Pueblo tribes, are said to have left their masks behind when sent to dwell in the bottom of a desert lake. Their return to help bring the rain is incarnated by the masked dancer. Cylindrical masks, covering the entire head and resting on the shoulders, are of a primal type. They are made of leather and humanized by the addition of hair and a variety of adjuncts. Eyes are represented by incisions or by buckskin balls filled with deer hair and affixed to the mask. The nose is often of rolled buckskin or corncob. Frequently the mask has a projecting wooden cylinder for a bill or a gourd stem cut with teeth for a snout. Horns are attached to some masks. Many colours are used in their painting; plumes and beads are attached, and the sex of the mask is distinguished by its shape: round head indicates male and square indicates female. In the western Sudan area of Africa many tribes have masked fertility ceremonials. The *segoni-kun* masks that are fashioned by the Bambara tribes in Mali are aesthetically among the most interesting. Antelopes, characterized by their elegant simplicity, are carved in wood and affixed to woven fibre caps that are hung with raffia and cover the wearer. The antelope is believed to have introduced agriculture, and so when crops are sown, members of Tjiwara society cavort in the fields in pairs to symbolize fertility and abundance.

Festive uses. Masks for festive occasions are still commonly used in the 20th century. Ludicrous, grotesque, or superficially horrible, festival masks are usually conducive to good-natured license, release from inhibitions, and ribaldry. These include the Halloween, Mardi Gras, or "masked ball" variety. The disguise is assumed to create a momentary, amusing character, often resulting in humorous confusions, or to achieve anonymity for the prankster or ribald reveller.

Throughout contemporary Europe and Latin America, masks are associated with folk festivals, especially those generated by seasonal changes or marking the beginning and end of the year. Among the most famous of the folk masks are the masks worn to symbolize the driving away of winter in parts of Austria and Switzerland. In Mexico and Guatemala, annual folk festivals employ masks for storytelling and caricature, such as for the Dance of the Old Men and the Dance of the Moors and the Christians. The Eskimo make masks with comic or satiric features

Uses in warfare and sport

Uses in fertility ceremonies



Agriculture fertility masks.

(Left) Corn husk mask of the Seneca Indians of the Iroquois nation of New York. Height 43 cm. (Right) Sekya, a *kachina* mask of the Zuni Indians of New Mexico. Painted leather, trimmed with feathers and hair. Height 43 cm.

By courtesy of the Museum of the American Indian, New York City



Detail from "The Ridotto" by Pietro Longhi (1702–85), showing the masks popularly worn at fashionable Venetian functions of the 18th century. In the Galleria dell'Accademia Carrara, Bergamo, Italy.

SCALA—Art Resource

that are worn at festivals of merrymaking, as do the *lbos* of Nigeria.

Theatrical uses. Masks have been used almost universally to represent characters in theatrical performances. Theatrical performances are a visual literature of a transient, momentary kind. It is most impressive because it can be seen as a reality; it expends itself by its very revelation. The mask participates as a more enduring element, since its form is physical.

The mask as a device for theatre first emerged in Western civilization from the religious practices of ancient Greece. In the worship of Dionysus, god of fecundity and the harvest, the communicants' attempt to impersonate the deity by donning goatskins and by imbibing wine eventually developed into the sophistication of masking. When a literature of worship appeared, a disguise, which consisted of a white linen mask hung over the face (a device supposedly initiated by Thespis, a 6th-century-BC poet who is credited with originating tragedy), enabled the leaders of the ceremony to make the god manifest. Thus symbolically identified, the communicant was inspired to speak in the first person, thereby giving birth to the art of drama. In Greece the progress from ritual to ritual-drama was continued in highly formalized the-

atrical representations. Masks used in these productions became elaborate headpieces made of leather or painted canvas and depicted an extensive variety of personalities, ages, ranks, and occupations. Heavily coiffured and of a size to enlarge the actor's presence, the Greek mask seems to have been designed to throw the voice by means of a built-in megaphone device and, by exaggeration of the features, to make clear at a distance the precise nature of the character. Moreover, their use made it possible for the Greek actors—who were limited by convention to three speakers for each tragedy—to impersonate a number of different characters during the play simply by changing masks and costumes. Details from frescoes, mosaics, vase paintings, and fragments of stone sculpture that have survived to the present day provide most of what is known of the appearance of these ancient theatrical masks. The tendency of the early Greek and Roman artists to idealize their subjects throws doubt, however, upon the accuracy of these reproductions. In fact, some authorities maintain that the masks of the ancient theatre were crude affairs with little aesthetic appeal.

In the Middle Ages, masks were used in the mystery plays of the 12th to the 16th century. In plays dramatizing portions of the Old and New Testaments, grotesques of all sorts, such as devils, demons, dragons, and personifications of the seven deadly sins, were brought to stage life by the use of masks. Constructed of papier-mâché, the masks of the mystery plays were evidently marvels of ingenuity and craftsmanship, being made to articulate and to belch fire and smoke from hidden contrivances. But again, no reliable pictorial record has survived. Masks used in connection with present-day carnivals and Mardi Gras and those of folk demons and characters still used by central European peasants, such as the *Perchten* masks of Alpine Austria, are most likely the inheritors of the tradition of medieval masks.

The 15th-century Renaissance in Italy witnessed the rise of a theatrical phenomenon that spread rapidly to France, to Germany, and to England, where it maintained its popularity into the 18th century. Comedies improvised from scenarios based upon the domestic dramas of the ancient Roman comic playwrights Plautus (254?–184 BC) and Terence (186/185–159 BC) and upon situations drawn from anonymous ancient Roman mimes flourished under the title of *commedia dell'arte*. Adopting the Roman stock figures and situations to their own usages, the players of the *commedia* were usually masked. Sometimes the masking was grotesque and fanciful, but generally a heavy leather mask, full or half face, disguised the *commedia* player. Excellent pictorial records of both *commedia* costumes and masks exist; some sketches show the characters of Arlecchino and Colombina wearing black masks covering merely the eyes, from which the later masquerade mask is certainly a development.

Except for vestiges of the *commedia* in the form of puppet and marionette shows, the drama of masks all but disappeared in Western theatre during the 18th, 19th, and first half of the 20th centuries. In modern revivals of ancient Greek plays, masks have occasionally been employed, and such highly symbolic plays as *Die versunkene Glocke* (*The Sunken Bell*; 1897) by the German Gerhart Hauptmann (1862–1946) and dramatizations of *Alice in Wonderland* have required masks for the performers of grotesque or animal figures. The Irish poet-playwright W.B. Yeats (1865–1939) revived the convention in his *Dreaming of the Bones* and in other plays patterned upon the Japanese *Nō* drama. In 1926 theatregoers in the United States witnessed a memorable use of masks in *The Great God Brown* by the American dramatist Eugene O'Neill (1888–1953), wherein actors wore masks of their own faces to indicate changes in the internal and external lives of their characters. Oskar Schlemmer (1888–1943), a German artist associated with the Bauhaus, became interested in the late 1920s and '30s in semantic phenomenology as applied to the design of masks for theatrical productions. Modern art movements are often reflected in the design of contemporary theatrical masks. The stylistic concepts of Cubism and Surrealism, for example, are apparent in the masks executed for a 1957 production of *La favola*

Ancient
Greek and
Roman

Medieval
and
Renaissance

Modern



Masked Roman actors in a comedy scene, 1st-century relief from Pompeii.

del figlio cambiato (*The Fable of the Transformed Son*) by the Italian dramatist Luigi Pirandello (1867–1936). A well-known mid-century play using masks was *Les Nègres* (1958; *The Blacks*, 1960) by the French writer Jean Genet. The mask, however, has unquestionably lost its importance as a theatrical convention in the 20th century, and its appearance in modern plays is unusual.

East Asian

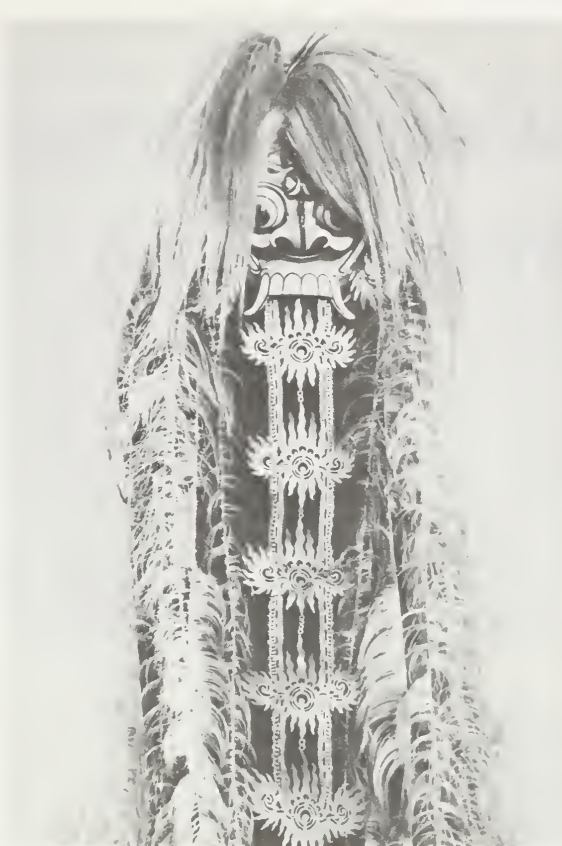
In many ways akin to Greek drama in origin and theme, the Nō drama of Japan has remained a significant part of national life since its beginnings in the 14th century. Nō masks, of which there are about 125 named varieties, are rigidly traditional and are classified into five general types: old persons (male and female), gods, goddesses, devils, and goblins. The material of the Nō mask is wood with a coating of plaster, which is lacquered and gilded. Colours are traditional. White is used to characterize a corrupt ruler; red signifies a righteous man; a black mask is worn by the villain, who epitomizes violence and brutality. Nō masks are highly stylized and generally characterized. They are exquisitely carved by highly respected artists known as *tenka-ichi*, “the first under heaven.” Shades of feeling are portrayed with beautifully sublimated realism. When the masks are subtly moved by the player’s hand or body motion, their expression appears to change.

In Tibet, sacred dramas are performed by masked lay actors. A play for exorcising demons called the “Dance of the Red Tiger Devil” is performed at fixed seasons of the year exclusively by the priests or lamas wearing awe-inspiring masks of deities and demons. Masks employed in this mystery play are made of papier-mâché, cloth, and occasionally gilt copper. In the Indian state of Sikkim and in Bhutan, where wood is abundant and the damp climate is destructive to paper, they are carved of durable wood. All masks of the Himalayan peoples are fantastically painted and are usually provided with wigs of yak tail in various colours. Formally they often emphasize the hideous.

Masks, usually made of papier-mâché, are employed in the religious or admonitory drama of China; but for the greater part the actors in popular or secular drama make up their faces with cosmetics and paint to resemble masks, as do the Kabuki actors in Japan. The makeup mask both identifies the particular character and conveys his personality. The highly didactic sacred drama of China is performed with the actors wearing fanciful and grotesque masks. Akin to this “morality” drama are the congratulatory playlets, pageants, processions, and dances of China. Masks employed in these ceremonies are highly ornamented, with jewelled and elaborately filigreed headgears. In the lion and dragon dances of both China and Japan, a stylized mask of the beast is carried on a pole by itinerant players, whose bodies are concealed by a dependent cloth. The mask and cloth are manipulated violently, as if the animal were in pursuit, to the taps of a small drum. The mask’s lower jaw is movable and made to emit a loud continuous clacking by means of a string.

Indonesian

On Java and Bali, wooden masks, *tupeng*, are used in certain theatrical performances called *wayang wong*. These dance dramas developed from the shadow puppet plays of the 18th century and are performed not only as amusement but as a safeguard against calamities. The stories are in part derived from ancient Sanskrit literature, especially the Hindu epics, although the Javanese later became Muslims. The brightly painted masks are made of wood and leather and are often fitted with horsehair and metallic or gilded paper accoutrements. They are ordinarily held in the teeth by means of a strap of leather or rattan that has been fastened across the inside. Occasionally an actor interrupts the unseen narrator, the Dalang, who is speaking the play. The mask is then held in front of the face while the player says his line. The use of theatrical masks in Java is exceptional, since masks, being forbidden under the prohibition of images, are practically unknown in the Islāmic world.



Javanese *tupeng* mask representing the witch Rangda. Lacquered wood, cloth, metal, and horsehair.

By courtesy of the Tropenmuseum, Amsterdam

In the 20th century, with the breaking down of primitive and folk cultures, the mask has increasingly become a decorative object, although it has long been used in art as an ornamental device. In Haiti, India, Indonesia, Japan, Kenya, and Mexico, masks are produced largely for tourists. The collecting of old masks has been a part of the current interest in so-called primitive and folk arts. Masks also have exerted a decided influence on modern art movements, especially in the first decades of the 20th century, when painters in France and Germany found a source of inspiration in the tribal masks of Africa and western Oceania.

BIBLIOGRAPHY. “Masks” in the *Encyclopedia of World Art*, vol. 9, col. 520–570 (1964), a good historical survey; WILLIAM N. FENTON, “Masked Medicine Societies of the Iroquois,” *Smithsonian Institution, Annual Report for 1940*, pp. 397–429 (1941), a very good general discussion of Iroquois masks, with illustrations; MARCEL GRIAULE, *Masques Dogons*, 2nd ed. (1963), a profusely illustrated classic study of the masks of the Dogon, a people of Mali, within their cultural setting; GEORGE W. HARLEY, *Masks as Agents of Social Control in Northeast Liberia*, Peabody Museum Papers 32, no. 2 (1950, reprinted 1975), a useful, illustrated article on this topic; EDWARD A. KENNARD, *Hopi Kachinas*, 2nd ed. (1971), an important study; DOROTHY J. RAY, *Eskimo Masks: Art and Ceremony* (1967), one of the best studies of Eskimo masks, with many fine plates and a bibliography; F.E. WILLIAMS, *Drama of Orokoloko* (1940, reprinted 1969), a classic study of masks of the Gulf of Papua, New Guinea, with fine illustrations and a bibliography; MALCOLM KIRK, *Man as Art: New Guinea* (1981), with especially good photographs; DONALD B. CORDRY, *Mexican Masks* (1980), a study of how Mexican masks are related to both the European and the Indian traditions; SIMON OTTENBERG, *Masked Rituals of Afikpo* (1975), a survey of a Nigerian masquerade tradition; and LEON UNDERWOOD, *Masks of West Africa* (1948), a small book but important for the subject, with good plates.

(P.S.W.)

Materials Science

Materials science is the study of the properties of solid materials and how those properties are determined by a material's composition and structure. It grew out of an amalgam of solid-state physics, metallurgy, and chemistry, since the rich variety of materials properties cannot be understood within the context of any single classical discipline. With a basic understanding of the origins of properties, materials can be selected or designed for an enormous variety of applications, ranging from structural steels to computer microchips. Materials science is therefore important to many engineering activities such as electronics, aerospace, telecommunications, information processing, nuclear power, and energy conversion.

This article approaches the subject of materials science through five major fields of application: energy, ground transportation, aerospace, computers and communications, and medicine. The discussions focus on the funda-

mental requirements of each field of application and on the abilities of various materials to meet those requirements.

The many materials studied and applied in materials science are usually divided into four categories: metals, polymers, semiconductors, and ceramics. The sources, processing, and fabrication of these materials are explained at length in three *Macropædia* articles: INDUSTRIES, EXTRACTION AND PROCESSING; INDUSTRIES, CHEMICAL PROCESS; and INDUSTRIAL GLASS AND CERAMICS. Atomic and molecular structures are discussed in CHEMICAL ELEMENTS and MATTER. The applications covered in this article are given broad coverage in ENERGY CONVERSION, TRANSPORTATION, ELECTRONICS, and MEDICINE.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, sections 10/37, 111, 121, 125, 126, 334, 336, 712, 721, 724, 725, 734, and 735, and the *Index*.

This article is divided into the following sections:

Materials for energy 552	Composites
Classification of energy-related materials	Materials for computers and communications 559
Applications of energy-related materials	Electronic materials
Materials for ground transportation 555	Photonic materials
Metals	Materials for medicine 562
Plastics and composites	General requirements of biomaterials
Ceramics	Polymer biomaterials
Materials for aerospace 557	Applications of biomaterials
Metals	Bibliography 565

Materials for energy

An industrially advanced society uses energy and materials in large amounts. Transportation, heating and cooling, industrial processes, communications—in fact, all the physical characteristics of modern life—depend on the flow and transformation of energy and materials through the techno-economic system. These two flows are inseparably intertwined and form the lifeblood of industrial society. The relationship of materials science to energy usage is pervasive and complex. At every stage of energy production, distribution, conversion, and utilization, materials play an essential role, and often special materials properties are needed. Remarkable growth in the understanding of the properties and structures of materials enables new materials, as well as improvements of old ones, to be developed on a scientific basis, thereby contributing to greater efficiency and lower costs.

CLASSIFICATION OF ENERGY-RELATED MATERIALS

Energy materials can be classified in a variety of ways. For example, they can be divided into materials that are passive or active. Those in the passive group do not take part in the actual energy-conversion process but act as containers, tools, or structures such as reactor vessels, pipelines, turbine blades, or oil drills. Active materials are those that take part directly in energy conversion—such as solar cells, batteries, catalysts, and superconducting magnets.

Another way of classifying energy materials is by their use in conventional, advanced, and possible future energy systems. In conventional energy systems such as fossil fuels, hydroelectric generation, and nuclear reactors, the materials problems are well understood and are usually associated with structural mechanical properties or long-standing chemical effects such as corrosion. Advanced energy systems are in the development stage and are in actual use in limited markets. These include oil from shale and tar sands, coal gasification and liquefaction, photovoltaics, geothermal energy, and wind power. Possible future energy systems are not yet commercially deployed to any

significant extent and require much more research before they can be used. These include hydrogen fuel and fast-breeder reactors, biomass conversion, and superconducting magnets for storing electricity.

Classifying energy materials as passive or active or in relation to conventional, advanced, or future energy systems is useful because it provides a picture of the nature and degree of urgency of the associated materials requirements. But the most illuminating framework for understanding the relation of energy to materials is in the materials properties that are essential for various energy applications. Because of its breadth and variety, such a framework is best shown by examples. Accordingly, some of the most important materials properties that are needed for specific energy applications are listed in Table 1. The table could be greatly enlarged, but it is sufficient to show that inexpensive and abundant energy depends on the availability of materials with special properties. Each entry in the table states a critical property for the given application. For each application, many materials properties are utilized, but the table lists only those that are obviously most important. In oil refining, for example, reaction vessels must have certain mechanical and thermal properties, but catalysis is the critical process.

APPLICATIONS OF ENERGY-RELATED MATERIALS

High-temperature materials. In order to extract useful work from a fuel, it must first be burned so as to bring some fluid (usually steam) to high temperatures. Thermodynamics indicates that the higher the temperature, the greater the efficiency of the conversion of heat to work; therefore, the development of materials for combustion chambers, pistons, valves, rotors, and turbine blades that can function at ever-higher temperatures is of critical importance. The first steam engines had an efficiency of less than 1 percent, while modern steam turbines achieve efficiencies of 35 percent or more. Part of this improvement has come from improved design and metalworking accuracy, but a large portion is the result of using improved high-temperature materials. The early engines

Passive
and active
materials

Table 1: Materials Properties for Energy Applications

application	materials properties					
	mechanical	thermal	chemical	electric	magnetic	optical
Heat engine	high-temperature strength		corrosion resistance			
Oil-well drill bits	hardness					
Oil refining			catalysis			
Electricity generator	wear resistance and high-temperature strength			insulation	electromagnetism	
Nuclear fuel	swelling control					
Nuclear steam generator	resistance to stress corrosion and cracking		corrosion resistance			
Nuclear pressure vessel	resistance to crack growth					
Nuclear waste disposal			encapsulation; barriers			
Coal liquefaction and gasification	erosion resistance		catalysis			
Fuel cell			catalysis; corrosion resistance	electrolysis		
Solar energy		heat absorption		photoelectricity		reflectance
Wind power	fatigue resistance					
Geothermal energy			corrosion resistance			
Batteries				high energy density		
Superconductor	ductility; strength			high current capacity	magnetic quenching	
Hydrogen fuel			hydrogen absorption capacity	electrolysis		
Conservation	light weight; strength	thermal insulation; high-temperature resistance	low-temperature catalysis	semiconductivity	magnetic efficiency	low transmission loss

were made of cast iron and then ordinary steels. Later, high-temperature alloys containing nickel, molybdenum, chromium, and silicon were developed that did not melt or fail at temperatures above 540° C (1,000° F). But modern combustion processes are nearing the useful temperature limits that can be achieved with metals, and so new materials that can function at higher temperatures—particularly intermetallic compounds and ceramics—are being developed.

The structural features that limit the use of metals at high temperatures are both atomic and electronic. All materials contain dislocations. The simplest of these are the result of planes of atoms that do not extend all through the crystal, so that there is a line where the plane ends that has fewer atoms than normal. In metals, the outer electrons are free to move. This gives a delocalized cohesion so that, when a stress is applied, dislocations can move to relieve the stress. The result is that metals are ductile: not only can they be easily worked into desired shapes, but when stressed they will gradually yield plastically rather than breaking immediately. This is a desirable feature, but the higher the temperature, the greater the plastic flow under stress—and, if the temperature is too high, the material will become useless. In order to get around this, materials are being studied in which the motion of dislocations is inhibited. Ceramics such as silicon nitride or silicon carbide and intermetallics such as nickel aluminide hold promise because the electrons that hold them together are highly localized in the form of valence or ionic bonds. It is as if metals were held together by a slippery glue while in nonmetals the atoms were connected by rigid rods. Dislocations thus find it much harder to move in nonmetals; raising the temperature does not increase dislocation motion, and the stress needed to make them yield is much higher. Furthermore, their melting points are significantly higher than those of metals, and they are much more resistant to chemical attack. But these desirable features come at a price. The very structure that makes them attractive also makes them brittle; that is, they do not flow

when subject to a high stress and are prone to failure by cracking. Modern research is aimed at overcoming this lack of ductility by modification of the material and how it is made. Hot pressing of ceramic powders, for example, minimizes the number of defects at which cracks can start, and the addition of small amounts of certain metals to intermetallics strengthens the cohesion among crystal grains at which fractures normally develop. Such advances, along with intelligent design, hold the promise of being able to build heat engines of much higher efficiency than those now available.

Diamond drills. Diamond drill bits are an excellent example of how an old material can be improved. Diamond is the hardest known substance and would make an excellent drill bit except that it is expensive and has weak planes in its crystal structure. Because natural diamonds are single crystals, the planes extend throughout the material, and they cleave easily. Such cleavage planes allow a diamond cutter to produce beautiful gems, but they are a disaster for drilling through rock. This limitation was overcome by Stratapax, a sintered diamond material developed by the General Electric Company of the United States. This consists of synthetic diamond powder that is formed into a thin plate and bonded to tungsten-carbide studs by sintering (fusing by heating the material below the melting point). Because the diamond plate is polycrystalline, cleavage cannot propagate through the material. The result is a very hard bit that does not fail by cleavage when it is used to drill through rock to get at oil and natural gas.

Oil platforms. An important example of dealing with old problems by modern methods is provided by the prevention of crack growth in offshore oil-drilling platforms. The primary structure consists of welded steel tubing that is subject to continually varying stress from ocean waves. Since the cost of building and deploying a platform can amount to several billion dollars, it is imperative that the platform have a long life and not be lost because of premature metal failure.

The problem of dislocations in metals

The growth of cracks

In the North Sea, 75 percent of the waves are higher than two metres (six feet) and exert considerable stresses on the platform. Cyclic loading of a metal ultimately results in fatigue failure in which surface cracks form, grow over time, and eventually cause the metal to break. Welds are the weak spots for such a process because weld metal has mechanical properties that are inferior to steel, and these are made even worse by internal stresses and defects (such as tiny voids and oxide particles) that are introduced in the welding process. Furthermore, the tube geometry at the weld consists of T- and K-shaped joints, which are natural stress concentrators. Fatigue failure in oil platforms therefore takes place at welds.

Fatigue occurs because cyclic stress causes dislocations to form and to move back and forth in the metal. Dislocation motion can be impeded by the presence of barriers such as small voids, grain boundaries, other dislocations, impurities, or even the surface itself. When dislocations are thereby pinned down, they stop the motion of other dislocations created by the stress, and a tangled dislocation network forms that results in a hard spot in the weld. The stress is then not easily relieved, and types of dislocation motion that are characteristic of the fatigue process initiate a crack at the weld surface. This phenomenon is a direct result of the microstructure of the weld and could be minimized by making the weld very uniform, preferably of the same material as the tubing, and having a very gently curved geometry at the joint. But, in spite of the sophistication of modern welding techniques, this is not yet feasible. An alternate strategy is therefore used in which the progress of the weld crack is monitored so that repairs can be made in time to avoid catastrophic failure. This can be done because, given the geometry of the joint, the depth of the crack is proportional to time until the crack is quite large. By contrast, in laboratory tests in which simple strips of metal are subject to cyclic stress, the growth rate increases as the crack becomes larger. In the T or K configuration in oil platforms, stress is much more evenly distributed, and the crack does not grow at an increasing speed until it is close to being fatal.

A technique for measuring the crack depth is based on the skin effect, the phenomenon in which a high-frequency alternating current is confined to the surface of a conductor. This makes it possible to measure the surface area of a small region with a simple meter, since an increase in crack depth means an increase in current path, and this in turn causes an increase in voltage drop. Measurement over time then allows the time to failure to be estimated; repairs can be effected before failure occurs. In this case, a knowledge of microstructure, the materials science of fatigue, and the study of crack formation have led to a simple testing technique of great economic importance.

Mathematical modeling of mass motion and heat transfer (including convection), along with studies of solidification, gas dissolution, and the effects of fluxes, are providing a much more detailed understanding of the factors controlling weld structure. With this knowledge, it should be possible to make welds with far fewer defects.

Radioactive waste. A different example is provided by the disposal of radioactive waste. Here the issue is primarily safety and the perception of safety rather than economics. Waste disposal will continue to be one of the factors that inhibit the exploitation of nuclear power until the public perceives it as posing no danger. The current plan is to interpose three barriers between the waste and human beings by first encapsulating it in a solid material, putting that in a metal container, and finally burying that container in geologically stable formations. The first step requires an inert, stable material that will hold the radioactive atoms trapped for a very long time, while the second step requires a material that is highly resistant to corrosion and degradation.

There are two good candidates for encapsulation. The first is borosilicate glass; this can be melted with the radioactive material, which then becomes a part of the glass structure. Glass has a very low solubility, and atoms in it have a very low rate of migration, so that it provides an excellent barrier to the escape of radioactivity. However, glass devitrifies at the high temperatures resulting from the

heat of radioactive decay; that is to say, the amorphous glassy state becomes crystalline, and, during this process, many cracks form in the material so that it no longer provides a good barrier against the escape of radioactive atoms. (This problem is more severe in rock than in salt formations, because salt has higher thermal conductivity than rock and dissipates the heat more easily.) The problem can be eased by storing the waste above ground for a decade or so. This would allow the initially high rate of decay to decrease, thereby lowering the temperature that would be reached after encapsulation. Handled in this way, borosilicate glass would be an excellent encapsulation material for reactor waste that had been aged for a decade or so.

The other candidate is a synthetic rock made of mineral mixtures such as zirconolite and perovskite. These are very insoluble and, in their natural state, are known to have sequestered radioactive elements for hundreds of millions of years. They are crystalline, ceramic materials whose crystal structures allow radioactive atoms to be immobilized within them. They are not subject to devitrification, since they are already crystalline.

Once encapsulated, radioactive waste must be put into canisters that are corrosion-resistant. These can be made of nickel-steel alloys, but the best candidate so far is a titanium material containing small amounts of nickel and molybdenum and traces of carbon and iron. Even though they are meant to be buried in as dry an environment as possible, these metals are tested by immersing them in brine. Tests show that seawater at 250° C (480° F) would corrode away less than one micrometre (one-thousandth of a millimetre, or four ten-thousandths of an inch) of the surface of the titanium material (known as Ti code 12) per year. This remarkable performance is primarily the result of a tough, highly resistant oxide skin that forms on titanium when exposed to oxygen. It would take thousands of years for the canisters to be penetrated by corrosion.

In order to estimate the effectiveness of such waste disposal, it must be noted that the waste is highly radioactive and dangerous initially but that the danger decreases with time. Radioactivity decays to such levels that the danger is much less after a few hundred years, extremely low after 500 years, and negligible after 1,000 years. In order to breach the triple-barrier system, groundwater must migrate to the canister, eat it away, and then leach out the radioactive atoms from the encapsulating glass or ceramic. This is a process that most probably would take far longer than a single millennium. A careful application of materials science can make radioactive waste disposal safer than current disposal methods for other toxic wastes.

Photovoltaics. Photovoltaic systems are an attractive alternative to fossil or nuclear fuels for the generation of electricity. Sunlight is free, it does not use up an irreplaceable resource, and its conversion to electricity is nonpolluting. In fact, photovoltaics are now in use where power lines from utility grids are either not possible or do not exist, as in outer space or remote, nonurban locations.

The barrier to widespread use of sunlight to generate electricity is the cost of photovoltaic systems. The application of materials science is essential in efforts to lower the cost to levels that can compete with those for fossil or nuclear fuels.

The conversion of light to electricity depends on the electronic structure of solar cells with two or more layers of semiconductor material that can absorb photons, the primary energy packets of light. The photons raise the energy level of the electrons in the semiconductor, exciting some to jump from the lower-energy valence band to the higher-energy conduction band. The electrons in the conduction band and the holes they have left behind in the valence band are both mobile and can be induced to move by a voltage. The electron motion, and the movement of holes in the opposite direction, constitute an electric current. The force that drives electrons and holes through a circuit is created by the junction of two dissimilar semiconducting materials, one of which has a tendency to give up electrons and acquire holes (thereby becoming the positive, or *p*-type, charge carrier) while the other accepts electrons (becoming the negative, or *n*-type, carrier). The

Encapsu-
lation

Radio-
active
decay

electronic structure that permits this is the band gap; it is equivalent to the energy required to move an electron from the lower band to the higher. The magnitude of this gap is important. Only photons with energy greater than that of the band gap can excite electrons from the valence band to the conduction band; therefore, the smaller the gap, the more efficiently light will be converted to electricity—since there is a greater range of light frequencies with sufficiently high energies. On the other hand, the gap cannot be too small, because the electrons and holes then find it easy to recombine, and a sizable current cannot be maintained.

Maximum efficiency of the solar cell

The band gap defines the theoretical maximum efficiency of a solar cell, but this cannot be attained because of other materials factors. For each material there is an intrinsic rate of recombination of electrons and holes that removes their contribution to electric current. This recombination is enhanced by surfaces, interfaces, and crystal defects such as grain boundaries, dislocations, and impurities. Also, a fraction of the light is reflected by the cell's surface rather than being absorbed, and some can pass through the cell without exciting electrons to the conduction band.

Improvements in the trade-off between cell efficiency and cost are well illustrated by the preparation of silicon that is the basic material of current solar cells. Initially, high-purity silicon was grown from a silicon melt by slowly pulling out a seed crystal that grew by the accretion and slow solidification of the molten material. Known as the Czochralski process, this resulted in a high-purity, single-crystal ingot that was then sliced into wafers about 1 millimetre (0.04 inch) thick. Each wafer's surface was then "doped" with impurities to create *p*-type and *n*-type materials with a junction between them. Metal was then deposited to provide electrical leads, and the wafer was encapsulated to yield a cell about 100 millimetres in diameter. This was an expensive and time-consuming process; it has been much improved in a variety of ways. For example, high-purity silicon can be made at drastically reduced cost by chemically converting ordinary silicon to silane or trichlorosilane and then reducing it back to silicon. This silane process is capable of continuous operation at a high production rate and with low energy input. In order to avoid the cost and waste associated with sawing silicon into wafers, methods of directly drawing molten silicon into thin sheets or ribbons have been developed; these can produce crystalline, polycrystalline, or amorphous material. Another alternative is the manufacture of thin films on ceramic substrates—a process that uses much less silicon than other methods. Single-crystal silicon has a higher efficiency than other forms, but it is also much more expensive. The materials challenge is to find a combination of cost and efficiency that makes photovoltaic electricity economically possible.

Surface treatments that increase efficiency include deposition of antireflecting coatings, such as silicon nitride, on the front of the cell and highly reflective coatings on the rear. Thus, more of the light that strikes a cell actually enters it, and light that escapes out the back is reflected back into the cell. An ingenious surface treatment is part of the point contact method, in which the surface of the cell is not planar but microgrooved so that light is randomly reflected as it strikes the cell. This increases the amount of light that can be captured by the cell. (L.A.G.)

Materials for ground transportation

The global effort to improve the efficiency of ground transportation vehicles, such as automobiles, buses, trucks, and trains, and thereby reduce the massive amounts of pollutants they emit, provides an excellent context within which to illustrate how materials science functions to develop new or better materials in response to critical human needs. For the automobile industry in particular, the story is a fascinating one in which the desire for lower vehicle weight, reduced emissions, and improved fuel economy has led to intense competition among aluminum, plastics, and steel companies for shares in the enormous markets involved (40 million to 50 million cars and trucks per year worldwide). In this battle, materials scientists have a key

Reducing weight, emissions, and fuel consumption

role to play because the success of their efforts to develop improved materials will determine the shape and viability of future automobiles.

Just how seriously suppliers to the industry view the need either to protect or to increase their share of these enormous markets is demonstrated by their establishing of special programs, consortia, or centres that are specifically designed to develop better alloys, plastics, or ceramics for automotive applications. For example, in the United States a program at the Aluminum Company of America (Alcoa) called the aluminum intensive vehicle (AIV), and a similar one at Reynolds Metals, were established to develop materials and processes for making automobile "space frames" consisting of aluminum-alloy rods and die-cast connectors joined by welding and adhesive bonding. Not to be outdone, another aluminum company, Alcan Aluminium Limited of Canada, in a program entitled aluminum structured vehicle technology (ASVT), began to investigate the construction of automobile unibodies from adhesively bonded aluminum sheet. The plastics industry, of course, has a powerful interest in replacing as many metal automobile components as possible, and in order to help bring this about a centre called D&S Plastics International was formed in the Detroit, Mich., area of the United States by three corporations. The specific aim of this centre was to develop materials and a process suitable for forming several connected panels or components (e.g., body panels and bumper fascias) simultaneously out of different types of plastics. The centrepiece of the operation was a 4,000-ton co-injection press that could lead to cost reductions as great as 50 percent and thereby make the use of plastics for automotive applications more attractive.

In programs such as these, and in many more carried out by vendors and within the automobile companies themselves, materials scientists with specialized training in advanced metals, plastics, and ceramics have been leading a revolution in the automotive industry. The following sections describe specific needs that have been identified for improving the performance of automobiles and other ground-transportation vehicles, as well as approaches that materials scientists have taken in response to those needs.

METALS

Aluminum. Since aluminum has about one-third the density of steel, its substitution for steel in automobiles would seem to be a sensible approach to reducing weight and thereby increasing fuel economy and reducing harmful emissions. Such substitutions cannot be made, however, without due consideration of significant differences in other properties of the two materials. This is one important facet of the materials scientist's job—to help evaluate the suitability of a material for a given application based on how its properties balance against load and performance requirements specified by the design engineer. In this case (aluminum versus steel), it is instructive to consider the materials scientist's approach to evaluating the use of aluminum in automotive panels—such components as doors, hoods, trunk decks, and roofs that can make up more than 60 percent of a vehicle's weight.

Two primary properties of any metal are (1) its yield strength, defined as its ability to resist permanent deformation (such as a fender dent), and (2) its elastic modulus, defined as its ability to resist elastic or springy deflection like a drum head. By alloying, aluminum can be made to have a yield strength equal to a moderately strong steel and therefore to exhibit similar resistance to denting in an automobile panel. On the other hand, alloying does not normally affect the elastic modulus of metals significantly, so that automotive door panels or hoods made from aluminum alloys, all of which have approximately one-third the modulus of steel, would be floppy and suffer large deflections when buffeted by the wind, for example. From this point of view, aluminum would appear to be a marginal choice for body panels.

One might attempt to overcome this deficiency by increasing the thickness of the aluminum sheet stock to three times the thickness of the steel it is intended to replace. This, however, would simply increase the weight to roughly that of an equivalent steel structure and thus

Maximizing yield strength and elastic modulus

defeat the purpose of the exercise. Fortunately, as was elegantly demonstrated in 1980 by two British materials scientists, Michael Ashby and David Jones, when proper account is taken of the way an actual door panel deflects, constrained as it is by the door edges, it is possible to use aluminum sheet only slightly thicker than the steel it would replace and still achieve equivalent performance. The net result would be a weight savings of almost two-thirds by the substitution of aluminum for steel on such body components. This suggests that understanding the interrelationship between materials properties and structural design is an important factor in the successful application of materials science.

Another important activity of the materials scientist is that of alloy development, which in some cases involves designing alloys for very specific applications. For example, in Alcoa's AIV effort, materials scientists and engineers developed a special casting alloy for use as cast aluminum nodes (connectors) in their space frame design. Ordinarily, metal castings exhibit very little toughness, or ductility, and they are therefore prone to brittle fracture followed by catastrophic failure. Since the integrity of an automobile would be limited by having relatively brittle body components, a proprietary casting alloy and processing procedure were developed that provide a material of much greater ductility than is normally available in a casting alloy.

Many other advances in aluminum technology, brought about by materials scientists and design engineers, have led to a greater acceptance of aluminum in automobiles, trucks, buses, and even light rail vehicles. Among these are alloys for air-conditioner components that are designed to be chemically compatible with environmentally safer refrigerants and to withstand the higher pressures required by them. Also, alloys have been developed that combine good formability and corrosion resistance with the ability to achieve maximum strength without heat treating; these alloys develop their strength during the forming operation.

As a consequence, the list of vehicles that contain significant quantities of aluminum substituted for steel has steadily grown. A milestone was reached in 1992 with a limited-edition Jaguar sports car that was virtually all aluminum, including the engine, adhesively bonded chassis, and skin. Somewhat less expensive and in full production were Honda's Acura NSX, containing more than 400 kilograms (900 pounds) of aluminum compared with about 70 kilograms for the average automobile, and General Motors' Saturn, with an aluminum engine block and cylinder heads. These vehicles and others took their place alongside the British Land Rover, which was built with all-aluminum body panels beginning in 1948—a choice dictated by a shortage of steel during World War II and continued by the manufacturer ever since.

Steel. While the goal of the aluminum and plastics industries is to achieve vehicle weight reductions by substituting their products for steel components, the goal of the steel industry is to counter such inroads with such innovative developments as high-strength, but inexpensive, "microalloyed" steels that achieve weight savings by thickness reductions. In addition, alloys have been developed that can be tempered (strengthened) in paint-baking ovens rather than in separate and expensive heat-treatment furnaces normally required for conventional steels.

The microalloyed steels, also known as high-strength low-alloy (HSLA) steels, are intermediate in composition between carbon steels, whose properties are controlled mainly by the amount of carbon they contain (usually less than 1 percent), and alloy steels, which derive their strength, toughness, and corrosion resistance primarily from other elements, including silicon, nickel, and manganese, added in somewhat larger amounts. Developed in the 1960s and resurrected in the late 1970s to satisfy the need for weight savings through greater strength, the HSLA steels tend to be low in carbon with minute additions of titanium or vanadium, for example. Offering tensile strengths that can be triple the value of the carbon steels they are designed to replace (e.g., 700 megapascals versus 200 megapascals), they have led to significant weight savings through thickness reductions—albeit at a slight loss of structural

stiffness, because their elastic moduli are the same as other steels. They are considered to be quite competitive with aluminum substitutes for two reasons: they are relatively inexpensive (steel sells for one-half the price of aluminum on a per-unit-weight basis); and very little change in fabrication and processing procedures is needed in switching from carbon steel to HSLA steel, whereas major changes are usually required in switching to aluminum.

Bake-hardenable steels were developed specifically for the purpose of eliminating an expensive fabrication step—i.e., the heat-treating furnace, where steels are imparted with their final strength. To do this, materials scientists have designed steels that can be strengthened in the same ovens used to bake body paint onto the part. These furnaces must operate at relatively low temperatures (170° C, or 340° F), so that special steels had to be developed that would achieve suitable strengths at heat-treatment temperatures very much below those normally employed (up to 600° C, or 1,100° F). Knowing that high-alloy steels would never be hardenable at such low temperatures, materials scientists focused their attention on carbon steels, but even here adequate strengths could not be obtained initially. Then in the 1980s scientists at the Japanese Sumitomo Metal Industries developed a steel containing nitrogen (a gas that constitutes three-quarters of the Earth's atmosphere) in addition to carbon and several other additives. Very high strengths (over 900 megapascals) and excellent toughness can be achieved on formed parts with this inexpensive addition after baking for 20 minutes at temperatures typical for a paint-baking operation.

PLASTICS AND COMPOSITES

The motive for replacing the metal components of cars, trucks, and trains with plastics is the expectation of large weight savings due to the large differences in density involved: plastics are one-sixth the weight of steel and one-half that of aluminum per unit volume. However, as in evaluating the suitability of replacing steel with aluminum, the materials scientist must compare other properties of the materials in order to determine whether the tradeoffs are reasonable. For two reasons, the likely conclusion would be that plastics simply are not suitable for this type of application: the strength of most plastics, such as epoxies and polyesters, is roughly one-fifth that of steel or aluminum; and their elastic modulus is one-sixtieth that of steel and one-twentieth that of aluminum. On this basis, plastics do not appear to be suitable for structural components. What, then, accounts for the successful use that has been made of them? The answer lies in efforts made over the years by materials scientists, polymer chemists, mechanical engineers, and production managers to combine relatively weak and low-stiffness resins with high-strength, high-modulus reinforcements, thereby making new materials called composites with much more suitable properties than plastics alone.

The reinforcements used in composites are generally chosen for their high strength and modulus, as might be expected, but economic considerations often force compromises. For example, carbon fibres have extremely high modulus values (up to five times that of steel) and therefore make excellent reinforcements. However, their cost precludes their extensive use in automobiles, trucks, and trains, although they are used regularly in the aerospace industry. More suitable for non-aerospace applications are glass fibres (whose modulus can approach 1.5 times that of aluminum) or, in somewhat special cases, a mixture of glass and carbon fibres.

The physical form and shape of the reinforcements vary greatly, depending on many factors. The most effective reinforcements are long fibres, which are employed either in the form of a woven cloth or as separate layers of unidirectional fibres stacked upon one another until the proper laminate thickness is achieved. The resin may be applied to the fibres or cloth before laying up, thus forming what are termed prepregs, or it may be added later by "wetting out" the fibres. In either case, the assembly is then cured, usually under pressure, to form the composite. This type of composite takes full advantage of the properties of the fibres and is therefore capable of yielding strong, stiff

Bake-hardenable steels

The growing use of aluminum

Long-fibre reinforcement

panels. Unfortunately, the labour involved in the lay-up operations and other factors make it very expensive, so that long-fibre reinforcement is used only sparingly in the automobile industry.

One attempt to avoid expensive hand lay-up operations involves chopped fibres that are employed in mat form, somewhat like felt, or as loose fibres that may be either blown into a mold or injected into a mold along with the resin. Another method does not use fibres at all; instead the reinforcement is in the form of small, high-modulus particles. These are the least expensive of all to process, since the particles are simply mixed into the resin, and the mixture is used in various types of molds. On the other hand, particles are the least efficient reinforcement material; as a consequence, property improvements are not outstanding.

In choosing the other major constituent in composites, the polymer matrix, one faces a somewhat daunting variety, including epoxies, polyimides, polyurethanes, and polyesters. Each has its advantages and disadvantages that must be evaluated in order to determine suitability for a particular application. Among the factors to be considered are cost, processing temperature (curing temperature if using a thermoset polymer and melting temperature if using a thermoplastic), flow properties in the molding operation, sag resistance during paint bake out, moisture resistance, and shelf life. The number of combinations of resins, reinforcements, production methods, and fibre-to-resin ratios is so challenging that materials scientists must join forces with polymer chemists and engineers from the design, production, and quality-control departments of the company in order to choose the right combination for the application.

Judging by the inroads that have been made in replacing metals with composites, it appears that technologists have been making the right choices. The introduction of fiberglass-reinforced plastic skins on General Motors' 1953 Corvette sports car marked the first appearance of composites in a production model, and composites have continued to appear in automotive components ever since. In 1984, General Motors' Fiero was placed on the market with the entire body made from composites, and the Camaro/Firebird models followed with doors, roof panels, fenders, and other parts made of composites. Composites were also chosen for exterior panels in the Saturn, which appeared in 1990. In addition, they have had less visible applications—for example, the glass-reinforced nylon air-intake manifold on some BMW models.

CERAMICS

Ceramics play an important role in engine efficiency and pollution abatement in automobiles and trucks. For example, one type of ceramic, cordierite (a magnesium aluminosilicate), is used as a substrate and support for catalysts in catalytic converters. It was chosen for this purpose because, along with many ceramics, it is lightweight, can operate at very high temperatures without melting, and conducts heat poorly (helping to retain exhaust heat for improved catalytic efficiency). In a novel application of ceramics, a cylinder wall was made of transparent sapphire (aluminum oxide) by General Motors' researchers in order to examine visually the internal workings of a gasoline engine combustion chamber. The intention was to arrive at improved understanding of combustion control, leading to greater efficiency of internal-combustion engines.

Another application of ceramics to automotive needs is a ceramic sensor that is used to measure the oxygen content of exhaust gases. The ceramic, usually zirconium oxide to which a small amount of yttrium has been added, has the property of producing a voltage whose magnitude depends on the partial pressure of oxygen surrounding the material. The electrical signal obtained from such a sensor is then used to control the fuel-to-air ratio in the engine in order to obtain the most efficient operation.

Because of their brittleness, ceramics have not been used as load-bearing components in ground-transportation vehicles to any great extent. The problem remains a challenge to be solved by materials scientists of the future.

(J.D.V.)

Materials for aerospace

The primary goal in the selection of materials for aerospace structures is the enhancement of fuel efficiency to increase the distance traveled and the payload delivered. This goal can be attained by developments on two fronts: increased engine efficiency through higher operating temperatures and reduced structural weight. In order to meet these needs, materials scientists look to materials in two broad areas—metal alloys and advanced composite materials. A key factor contributing to the advancement of these new materials is the growing ability to tailor materials to achieve specific properties.

Enhancing
fuel
efficiency

METALS

Many of the advanced metals currently in use in aircraft were designed specifically for applications in gas-turbine engines, the components of which are exposed to high temperatures, corrosive gases, vibration, and high mechanical loads. During the period of early jet engines (from about 1940 to 1970), design requirements were met by the development of new alloys alone. But the more severe requirements of advanced propulsion systems have driven the development of novel alloys that can withstand temperatures greater than 1,000° C (1,800° F), and the structural performance of such alloys has been improved by developments in the processes of melting and solidification.

Melting and solidifying. Alloys are substances composed of two or more metals or of a metal and a nonmetal that are intimately united, usually by dissolving in each other when they are melted. The principal objectives of melting are to remove impurities and to mix the alloying ingredients homogeneously in the base metal. Major advances have been made with the development of new processes based on melting under vacuum (hot isostatic pressing), rapid solidification, and directional solidification.

In hot isostatic pressing, prealloyed powders are packed into a thin-walled, collapsible container, which is placed in a high-temperature vacuum to remove adsorbed gas molecules. It is then sealed and put in a press, where it is exposed to very high temperatures and pressures. The mold collapses and welds the powder together in the desired shape.

Molten metals cooled at rates as high as a million degrees per second tend to solidify into a relatively homogeneous microstructure, since there is insufficient time for crystalline grains to nucleate and grow. Such homogeneous materials tend to be stronger than the typical "grainy" metals. Rapid cooling rates can be achieved by "splat" cooling, in which molten droplets are projected onto a cold surface. Rapid heating and solidification can also be achieved by passing high-power laser beams over the material's surface.

Unlike composite materials (see below), grainy metals exhibit properties that are essentially the same in all directions, so they cannot be tailored to match anticipated load paths (*i.e.*, stresses applied in specific directions). However, a technique called directional solidification provides a certain degree of tailorability. In this process the temperature of the mold is precisely controlled to promote the formation of aligned stiff crystals as the molten metal cools. These serve to reinforce the component in the direction of alignment in the same fashion as fibres reinforce composite materials.

Alloying. These advances in processing have been accompanied by the development of new "superalloys." Superalloys are high-strength, often complex alloys that are resistant to high temperatures and severe mechanical stress and that exhibit high surface stability. They are commonly classified into three major categories: nickel-based, cobalt-based, and iron-based. Nickel-based superalloys predominate in the turbine section of jet engines. Although they have little inherent resistance to oxidation at high temperatures, they gain desirable properties through the addition of cobalt, chromium, tungsten, molybdenum, titanium, aluminum, and niobium.

Superalloys

Aluminum-lithium alloys are stiffer and less dense than conventional aluminum alloys. They are also "superplas-

Ceramics
in catalytic
converters

tic," owing to the fine grain size that can now be achieved in processing. Alloys in this group are appropriate for use in engine components exposed to intermediate to high temperatures; they can also be used in wing and body skins.

Titanium alloys, as modified to withstand high temperatures, are seeing increased use in turbine engines. They are also employed in airframes, primarily for military aircraft but to some extent for commercial planes as well.

COMPOSITES

While developments in metals have had an impact on engine design, there is a growing trend toward the application of composite materials to aerospace structures. One of the reasons for this is that alloys do not offer substantial weight savings, which is a primary advantage of composites. Indeed, advanced composites have been used most widely where saving mass results in either significantly improved performance or significantly lower life-cycle costs. The most extensive application, therefore, has been in satellite systems, military aircraft, radomes, helicopters, commercial transport aircraft, and general aviation.

Broadly defined, composites are materials with two or more distinct components that combine to yield characteristics superior to those of the individual constituents. Although this definition can apply to such ordinary building materials as plywood, concrete, and bricks, within the aerospace industry the term composite generally refers to the fibre-reinforced metal, polymer, and ceramic products that have come into use since World War II. These materials consist of fibres (such as glass, graphite, silicon carbide, or aramid) that are embedded in a matrix of, for example, aluminum, epoxy, or silicon nitride.

In the late 1950s a revolution in materials development occurred in response to the space program's need for lightweight, thermally stable materials. Boron-tungsten filaments, carbon-graphite fibres, and organic aramid fibres proved to be strong, stiff, and light, but one problem with using them as fibres was that they were of limited value in any construction other than rope, which can bear loads in only one direction. Materials scientists needed to develop a way to make them useful under all loading conditions, and this led to the development of composites. While the structural value of a bundle of fibres is low, the strength of individual fibres can be harnessed if they are embedded in a matrix that acts as an adhesive, binding the fibres and lending solidity to the material. The matrix also protects the fibres from environmental stress and physical damage, which can initiate cracks. In addition, while the strength and stiffness of the composite remain largely a function of the reinforcing material—that is, the fibres—the matrix can contribute other properties, such as thermal and electrical conductivity and, most important, thermal stability. Finally, fibre-matrix combination reduces the potential for complete fracture. In a monolithic (or single) material, a crack, once started, generally continues to propagate until the material fails; in a composite, if one fibre in an assemblage fails, the crack may not extend to the other fibres, so the damage is limited.

To some extent, the composite-materials engineer is trying to mimic structures made spontaneously by plants and animals. A tree, for example, is made of a fibre-reinforced material whose strength is derived from cellulose fibres that grow in directions that match the weight of the branches. Similarly, many organisms naturally fabricate "bioceramics," such as those found in shells, teeth, and bones. While the designers of composites for the aerospace industry would like to copy some of the features of bioceramics production—room-temperature processing and net-shape products, for example—they do not want to be constrained by slow processing methods and limited fibre and matrix material choices. In addition, unlike a mollusk, which has to produce only one shell, the composites manufacturer has to use rapid, repeatable processing methods that can fabricate hundreds or even thousands of parts.

Modern composites are generally classified into three categories according to the matrix material: polymer, metal, or ceramic. Since polymeric materials tend to degrade at elevated temperatures, polymer-matrix composites (PMCs)

are restricted to secondary structures in which operating temperatures are lower than 300° C (570° F). For higher temperatures, metal-matrix and ceramic-matrix composites are required.

Polymer-matrix composites. PMCs are of two broad types, thermosets and thermoplastics. Thermosets are solidified by irreversible chemical reactions, in which the molecules in the polymer "cross-link," or form connected chains. The most common thermosetting matrix materials for high-performance composites used in the aerospace industry are the epoxies. Thermoplastics, on the other hand, are melted and then solidified, a process that can be repeated numerous times for reprocessing. Although the manufacturing technologies for thermoplastics are generally not as well developed as those for thermosets, thermoplastics offer several advantages. First, they do not have the shelf-life problem associated with thermosets, which require freezer storage to halt the irreversible curing process that begins at room temperature. Second, they are more desirable from an environmental point of view, as they can be recycled. They also exhibit higher fracture toughness and better resistance to solvent attack. Unfortunately, thermoplastics are more expensive, and they generally do not resist heat as well as thermosets; however, strides are being made in developing thermoplastics with higher melting temperatures. Overall, thermoplastics offer a greater choice of processing approaches, so that the process can be determined by the scale and rate of production required and by the size of the component.

A variety of reinforcements can be used with both thermoset and thermoplastic PMCs, including particles, whiskers (very fine single crystals), discontinuous (short) fibres, continuous fibres, and textile preforms (made by braiding, weaving, or knitting fibres together in specified designs). Continuous fibres are more efficient at resisting loads than are short ones, but it is more difficult to fabricate complex shapes from materials containing continuous fibres than from short-fibre or particle-reinforced materials. To aid in processing, most high-performance composites are strengthened with filaments that are bundled into yarns. Each yarn, or tow, contains thousands of filaments, each of which has a diameter of approximately 10 micrometres (0.01 millimetre, or 0.0004 inch).

Depending on the application and on the type of load to be applied to the composite part, the reinforcement can be random, unidirectional (aligned in a single direction), or multidirectional (oriented in two or three dimensions). If the load is uniaxial, the fibres are all aligned in the load direction to gain maximum benefit of their stiffness and strength. However, for multidirectional loading (for example, in aircraft skins), the fibres must be oriented in a variety of directions. This is often accomplished by stacking layers (or lamina) of continuous-fibre systems.

The most common form of material used for the fabrication of composite structures is the prepregged tape, or "prepreg." There are two categories of prepreg: tapes, generally 75 millimetres (3 inches) or less in width, intended for fabrication in automated, computer-controlled tape-laying machines; and "broad goods," usually several metres in dimension, intended for hand lay-up and large sheet applications. To make prepregs, fibres are subjected to a surface treatment so that the resin will adhere to them. They are then placed in a resin bath and rolled into tapes or sheets.

To fabricate the composite, the manufacturer "lays up" the prepreg according to the reinforcement needs of the application. This has traditionally been done by hand, with successive layers of a broad-goods laminate stacked over a tool in the shape of the desired part in such a way as to accommodate the anticipated loads. However, efforts are now being directed toward automated fibre-placement methods in order to reduce costs and ensure quality and repeatability. Automated fibre-placement processes fall into two categories, tape laying and filament winding. The tape-laying process involves the use of devices that control the placement of narrow prepreg tapes over tooling with the contours of the desired part and along paths prescribed by the design requirements of the structure. The width of the tape determines the "sharpness" of the turns required

Thermosets
and
thermo-
plastics

Contribu-
tions of
fibre
reinforce-
ment
and the
matrix

Fabri-
cating a
composite
structure

to place the fibres in the prescribed direction—*i.e.*, wide tapes are used for gradual turns, while narrow tapes are required for the sharp turns associated with more complex shapes.

Filament winding uses the narrowest prepreg unit available—the yarn, or tow, of impregnated filaments. In this process, the tows are wound in prescribed directions over a rotating mandrel in the shape of the part. Successive layers are added until the required thickness is reached. Although filament winding was initially limited to geodesic paths (*i.e.*, winding the fibres along the most direct route between two points), the process is now capable of fabricating complex shapes through the use of robots.

For thermosetting polymers, the structure generated by either tape laying or filament winding must undergo a second manipulation in order to solidify the polymer through a curing reaction. This is usually accomplished by heating the completed structure in an autoclave, or oven. Thermoplastic systems offer the advantage of on-line consolidation, so that the high energy and capital costs associated with the curing step can be eliminated. For these systems, prepreg can be locally melted, consolidated, and cooled at the point of contact so that a finished structure is produced. A variety of energy sources are used to concentrate heat at the point of contact, including hot-gas torches, infrared light, and laser beams.

Pultrusion, the only truly continuous process for manufacturing parts from PMCs, is economical but limited to the production of beamlike shapes. On a pultrusion line, fibres and the resin are pushed through a heated die, or shaping tool, at one end, then cooled and pulled out at the other end. This process can be applied to both thermoplastic and thermoset polymers.

Resin transfer molding, or RTM, is a composites processing method that offers a high potential for tailorability but is currently limited to low-viscosity (easily flowing) thermosetting polymers. In RTM, a textile preform—made by braiding, weaving, or knitting fibres together in a specified design—is placed into a mold, which is then closed and injected with a resin. After consolidation, the mold is opened and the part removed. Preforms can be made in a wide variety of architectures, and several can be joined together during the RTM process to form a multi-element preform offering reinforcement in specific areas and load directions.

The similarity of meltable thermoplastic polymers to metals has prompted the extension of techniques used in metalworking. Sheet forming, used since the 19th century by metallurgists, is now applied to the processing of thermoplastic composites. In a typical thermoforming process, the sheet stock, or preform, is heated in an oven. At the forming temperature, the sheet is transferred into a forming system, where it is forced to conform to a tool, with a shape that matches the finished part. After forming, the sheet is cooled under pressure and then removed. Stretch forming, a variation on thermoplastic sheet forming, is specifically designed to take advantage of the extensibility, or ability to be stretched, of thermoplastics reinforced with long, discontinuous fibres. In this process, a straight preconsolidated beam is heated and then stretched over a shaped tool to introduce curvature. The specific advantage of stretch forming is that it provides an automated way to achieve a very high degree of fibre-orientation control in a wide range of part sizes.

Metal-matrix and ceramic-matrix composites. The requirement that finished parts be able to operate at temperatures high enough to melt or degrade a polymer matrix creates the need for other types of matrix materials, often metals. Metal matrices offer not only high-temperature resistance but also strength and ductility, or “bendability,” which increases toughness. The main problems with metal-matrix composites (MMCs) are that even the lightest metals are heavier than polymers, and they are very complex to process. MMCs can be used in such areas as the skin of a hypersonic aircraft, but on wing edges and in engines temperatures often exceed the melting point of metals. For the latter applications, ceramic-matrix composites (CMCs) are seeing increasing use, although the technology for CMCs is less mature than that for PMCs.

Ceramics consist of alumina, silica, zirconia, and other elements refined from fine earth and sand or of synthetic materials, such as silicon nitride or silicon carbide. The desirable properties of ceramics include superior heat resistance and low abrasive and corrosive properties. Their primary drawback is brittleness, which can be reduced by reinforcing with fibres or whiskers. The reinforcement material can be a metal or another ceramic.

Unlike polymers and metals, which can be processed by techniques that involve melting (or softening) followed by solidification, high-temperature ceramics cannot be melted. They are generally produced by some variation of sintering, a technique that renders a combination of materials into a coherent mass by heating to high temperatures without complete melting. If continuous fibres or textile weaves (as opposed to short fibres or whiskers) are involved, sintering is preceded by impregnating the assembly of fibres with a slurry of ceramic particles dispersed in a liquid. A major benefit of using CMCs in aircraft engines is that they allow higher operating temperatures and thus greater combustion efficiency, leading to reduced fuel consumption. An additional benefit is derived from the low density of CMCs, which translates into substantial weight savings.

Other advanced composites. Carbon-carbon composites are closely related to CMCs but differ in the methods by which they are produced. Carbon-carbon composites consist of semicrystalline carbon fibres embedded in a matrix of amorphous carbon. The composite begins as a PMC, with semicrystalline carbon fibres impregnated with a polymeric phenolic resin. The resin-soaked system is heated in an inert atmosphere to pyrolyze, or char, the polymer to a carbon residue. The composite is re-impregnated with polymer, and the pyrolysis is repeated. Continued repetition of this impregnation/pyrolysis process yields a structure with minimal voids. Carbon-carbon composites retain their strength at 2,500° C (4,500° F) and are used in the nose cones of reentry vehicles. However, because they are vulnerable to oxidation at such high temperatures, they must be protected by a thin layer of ceramic.

While materials research for aerospace applications has focused largely on mechanical properties such as stiffness and strength, other attributes are important for use in space. Materials are needed with a near-zero coefficient of thermal expansion; in other words, they have to be thermally stable and should not expand and contract when exposed to extreme changes in temperature. A great deal of research is focused on developing such materials for high-speed civilian aircraft, where thermal cycling is a major issue. High-toughness materials and nonflammable resin composite systems are also under investigation to improve the safety of aircraft interiors.

Efforts are also being directed toward the development of “smart,” or responsive, materials. Representing another attempt to mimic certain characteristics of living organisms, smart materials, with their built-in sensors and actuators, would react to their external environment by bringing on a desired response. This would be done by linking the mechanical, electrical, and magnetic properties of these materials. For example, piezoelectric materials generate an electrical current when they are bent; conversely, when an electrical current is passed through these materials, they stiffen. This property can be used to suppress vibration: the electrical current generated during vibration could be detected, amplified, and sent back, causing the material to stiffen and stop vibrating. (R.L.McC./D.S.K.)

Materials for computers and communications

The basic function of computers and communications systems is to process and transmit information in the form of signals representing data, speech, sound, documents, and visual images. These signals are created, transmitted, and processed as moving electrons or photons, and so the basic materials groups involved are classified as electronic and photonic. In some cases, materials known as optoelectronic bridge these two classes, combining abilities to interact usefully with both electrons and photons.

Among the electronic materials are various crystalline

Carbon-carbon composites

Forming of thermoplastic sheets

semiconductors; metalized film conductors; dielectric films; solders; ceramics and polymers formed into substrates on which circuits are assembled or printed; and gold or copper wiring and cabling.

Photonic materials include a number of compound semiconductors designed for light emission or detection; elemental dopants that serve as photonic performance-control agents; metal- or diamond-film heat sinks; metalized films for contacts, physical barriers, and bonding; and silica glass, ceramics, and rare earths for optical fibres.

ELECTRONIC MATERIALS

Between 1955 and 1990, improvements and innovations in semiconductor technology increased the performance and decreased the cost of electronic materials and devices by a factor of one million—an achievement unparalleled in the history of any technology. Along with this extraordinary explosion of technology has come an exponentially upward spiral of the capital investment necessary for manufacturing operations. In order to maintain cost-effectiveness and flexibility, radical changes in materials and manufacturing operations will be necessary.

Semiconductor crystals. *Silicon.* Bulk semiconductor silicon for the manufacture of integrated circuits (sometimes referred to as electronic-grade silicon) is the purest material ever made commercially in large quantities. One of the most important factors in preparing this material is control of such impurities as boron, phosphorus, and carbon (not to be confused with the dopants added later during circuit production). For the ultimate levels of integrated-circuit design, stray contaminant atoms must constitute less than 0.1 part per trillion of the material.

For fabrication into integrated circuits, bulk semiconductor silicon must be in the form of a single-crystal material with high crystalline perfection and the desired charge-carrier concentration. The size of the silicon ingot, or boule, has been scaled up in recent years, in order to provide wafers of increasing diameter that are demanded by the economics of integrated-circuit manufacturing. Most commonly, a 60-kilogram (130-pound) charge is grown to an ingot with a diameter of 200 millimetres (8 inches), but the semiconductor industry will soon require ingots as large as 300 millimetres. The ingots are then converted into wafers by machining and chemical processes.

III-V compounds. Although silicon is by far the most commonly used crystal material for integrated circuits, a significant volume of semiconductor devices and circuits employs III-V technology, so named because it is based on crystalline compounds formed by combining metallic elements from column III and nonmetallic elements from column V of the periodic table of chemical elements. When the elements are gallium and arsenic, the semiconductor is called gallium arsenide, or GaAs. However, other elements such as indium, phosphorus, and aluminum are often used in the compound to achieve specific performance characteristics.

For electronic applications, the III-V semiconductors offer the basic advantage of higher electron mobility, which translates into higher operating speeds. In addition, devices made with III-V compounds provide lower voltage operation for specific functions, radiation hardness (especially important for satellites and space vehicles), and semi-insulating substrates (avoiding the presence of parasitic capacitance in switching devices).

III-V materials are more difficult to handle than silicon, and a III-V wafer or substrate usually is less than half the size of a silicon wafer. In addition, a gallium arsenide wafer entering the processing facility can be expected to cost 10 to 20 times as much as a silicon wafer, although that cost difference narrows somewhat after fabrication, packaging, and testing. Nevertheless, there is one major characteristic of III-V materials with which silicon cannot compete: a III-V compound can be tailored to generate or detect photons of a specific wavelength. For example, an indium gallium arsenide phosphide (InGaAsP) laser can generate radiation at 1.55 micrometres to carry digitally coded information streams. (See below *Photonic materials*.) This means that a III-V component can fill both electronic and photonic functions in the same integrated circuit.

Photoresist films. Patterning polished wafers with an integrated circuit requires the use of photoresist materials that form thin coatings on the wafer before each step of the photolithographic process. Modern photoresists are polymeric materials that are modified when exposed to radiation (either in the form of visible, ultraviolet, or X-ray photons or in the form of energetic electron beams). A photoresist typically contains a photoactive compound (PAC) and an alkaline-soluble resin. The PAC, mixed into the resin, renders it insoluble. This mixture is coated onto the semiconductor wafer and is then exposed to radiation through a "mask" that carries the desired pattern. Exposed PAC is converted into an acid that renders the resin soluble, so that the resist can be dissolved and the exposed substrate beneath it chemically etched or metallically coated to match the circuit design.

Besides practical properties such as shelf life, cost, and availability, the key properties of a photoresist include purity, etching resistance, resolution, contrast, and sensitivity. As the feature sizes of integrated circuits shrink in each successive generation of microchips, photoresist materials are challenged to handle shorter wavelengths of light. For example, the photolithography of current designs (with features that have shrunk to less than one micrometre) is based on ultraviolet radiation in the wavelength range of 365 to 436 nanometres, but, in order to define accurately the smaller features of future microchips (less than 0.25 micrometre), shorter wavelengths will be necessary. The problem here is that electromagnetic radiation in such frequency regions is weaker. One solution is to use the chemically amplified photoresist, or CAMP. The sensitivity of a photoresist is measured by its quantum efficiency, or the number of chemical events that occur when a photon is absorbed by the material. In CAMP material, the number of events is dramatically increased by subsequent chemical reactions (hence the amplification), which means that less light is needed to complete the process.

Electric connections. The performance of today's electronic systems (and photonic systems as well) is limited significantly by interconnection technology, in which components and subsystems are linked by conductors and connectors. Currently, very fine gold or copper wiring, as thin as 30 micrometres, is used to carry electric current to and from the many pads along the sides or ends of a microchip to other components on a circuit board. The capacitance involved in such circuitry slows down the flow of electrons and, hence, of information. However, by integrating several chips into a single multichip module, in which the chips are connected on a shared substrate by various conducting materials (such as metalized film), the speed of information flow can be increased, thus improving the assembly's performance. Ideally, all the chips in a single module would be fabricated simultaneously on the same wafer, but in practice this is not feasible: Silicon crystal manufacture is still subject to an average of one flaw per wafer, meaning that at least one of the many chips cut from each wafer is scrapped. If the whole wafer area were dedicated to a single multifunction assembly, that one flaw would scrap the entire module. Multichip modules are therefore made up of as many as five microchips bonded to a silicon or ceramic substrate on which resistors and capacitors have been constructed with thin films. Typical materials used in a multichip module include the substrate; gold paste conductors applied in an additive process resembling silk screen printing; vitreous glazes to insulate the gold paste conductors from subsequent film layers; a series of thin films made with tantalum nitride, titanium, palladium, and plated gold; and a final package of silicone rubber.

Packaging materials. Several major types of packaging material are used by the electronics industry, including ceramic, refractory glass, premolded plastic, and postmolded plastic. Ceramic and glass packages cost more than plastic packages, so they make up less than 10 percent of the worldwide total. However, they provide the best protection for complex chips. Premolded plastic packages account for only a small but important fraction of the market, since they are required for packaging devices with many leads. Most plastic packages are postmolded, meaning that the

Control of
impurities

Photo-
active
compound

Multichip
modules

package body is molded over the assembly after the microchip has been attached to the fan-out pattern.

Precursors. The starting materials for most semiconductor devices are volatile and ultrapure gaseous derivatives of various organic and inorganic precursors. Many of them are toxic, and many will ignite spontaneously in the atmosphere. These gases are transported in high-pressure cylinders from the plant where they were made to the site where they will be used. One possible method of replacing these precursors with materials that are environmentally safe is known as *in situ* synthesis. In this method, dangerous reagents would be generated on demand in only the desired quantities, instead of being shipped cross-country and stored until needed at the semiconductor processing plant.

PHOTONIC MATERIALS

Computers and communications systems have been dominated by electronic technology since their beginnings, but photonic technology is making serious inroads throughout the information movement and management systems with such devices as lasers, light-emitting diodes, photodetecting diodes, optical switches, optical amplifiers, optical modulators, and optical fibres. Indeed, for long-distance terrestrial and transoceanic transmission of information, photonics has almost completely displaced electronics.

Crystalline materials. The light detectors and generators listed above are actually optoelectronic, because they link photonic and electronic systems. They employ the III-V compound semiconductors described above, many of them characterized by their band gaps—*i.e.*, the energy minimum of the electron conduction band and the energy maximum of hole valence bands occur at the same location in the momentum space, allowing electrons and holes to recombine and radiate photons efficiently. (By contrast, the conduction band minimum and the valence band maximum in silicon have dissimilar momenta, and therefore the electrons and holes cannot recombine efficiently.) Among the important compounds are gallium arsenide, aluminum gallium arsenide, indium gallium arsenide phosphide, indium phosphide, and aluminum indium arsenide.

Fabricating a single crystal from these combinations of elements is far more difficult than creating a single crystal of electronic-grade silicon. Special furnaces are required, and the process can take several days. Notwithstanding the precision involved, the sausage-shaped boule is less than half the diameter of a silicon ingot and is subject to a much higher rate of defects. Researchers are continuously seeking ways to reduce the thermal stresses that are primarily responsible for dislocations in the III-V crystal lattice that cause these defects. The purity and structural perfection of the final single-crystal substrates affect the qualities of the crystalline layers that are grown on them and the regions that are diffused or implanted in them during the manufacture of photonic devices.

Epitaxial layers. For the efficient emission or detection of photons, it is often necessary to constrain these processes to very thin semiconductor layers. These thin layers, grown atop bulk semiconductor wafers, are called epitaxial layers because their crystallinity matches that of the substrate even though the composition of the materials may differ—*e.g.*, gallium aluminum arsenide (GaAlAs) grown atop a gallium arsenide substrate. The resulting layers form what is called a heterostructure. Most continuously operating semiconductor lasers consist of heterostructures, a simple example consisting of 1000-angstrom thick gallium arsenide layers sandwiched between somewhat thicker (about 10000 angstroms) layers of gallium aluminum arsenide—all grown epitaxially on a gallium arsenide substrate. The sandwiching and repeating of very thin layers of a semiconductor between layers of a different composition allow one to modify the band gap of the sandwiched layer. This technique, called band-gap engineering, permits the creation of semiconductor materials with properties that cannot be found in nature. Band-gap engineering, used extensively with III-V compound semiconductors, can also be applied to elemental semiconductors such as silicon and germanium.

The most precise method of growing epitaxial layers on a semiconducting substrate is molecular-beam epitaxy (MBE). In this technique, a stream or beam of atoms or molecules is effused from a common source and travels across a vacuum to strike a heated crystal surface, forming a layer that has the same crystal structure as the substrate. Variations of MBE include elemental-source MBE, hydride-source MBE, gas-source MBE, and metal-organic MBE. Other approaches to epitaxial growth are liquid-phase epitaxy (LPE) or chemical vapour deposition (CVD). The latter method includes hydride CVD, trichloride CVD, and metal-organic CVD.

Normally, epitaxial layers are grown on flat surfaces, but scientists are searching for an economical and reliable method of growing epitaxial material on nonplanar structures—for example, around the “mesas” or “ridges” or in the “tubs” or “channels” that are etched into the surface of semiconducting devices. Nonplanar epitaxy is considered necessary for producing monolithic integrated optical devices or all-photonic switches and logic elements, but mastery of this method requires better understanding of the surface chemistry and surface dynamics of epitaxial growth.

Optical switching. Research in this area is driven by the need to switch data streams of higher and higher speed efficiently as customers for computer and communications services demand transmission and switching rates far higher than can be provided by a purely electronic system. Thanks to developments in semiconductor lasers and detectors (described above) and in optical fibres (described below), transmission at the desired high speeds has become possible. However, the switching of optical data streams still requires converting the data from the optical to the electronic domain, subjecting them to electronic switching and to manipulation inside the switching apparatus, and then reconverting the switched and reconfigured data into the optical domain for transmission over optical fibres. Electronic switching therefore is seen as the principal barrier to achieving higher switching speeds. One approach to solving this problem would be to introduce optics inside digital switching machines. Known as free-space photonics, this approach would involve such devices as semiconductor lasers or light-emitting diodes (LEDs), optical modulators, and photodetectors—all of which would be integrated into systems combined with electronic components.

One commercially available device for photonic switching is the quantum-well self-electro-optic-effect device, or SEED. The key concept for this device is the use of quantum wells. These structures consist of many thin layers of two different semiconductor materials. Individual layers are typically 10 nanometres (about 40 atoms) thick, and 100 layers are used in a device about 1 micrometre thick. When a voltage is applied across the layers, the transmission of photons through the quantum wells changes significantly, in effect creating an optical modulator—an essential component of any photonic circuit. Variations on the SEED concept are the symmetric SEED (S-SFED) and the field-effect transistor SEED. Neighbouring S-SEEDs could be connected by pairs of back-to-back quantum-well photodiodes, and commercially sized interconnection networks could be built by using free-space photonic interconnections between two-dimensional arrays of switching nodes. However, even this type of free-space optical interconnection technology would only enhance and extend electronic technology, not replace it.

The move of optoelectronic and photonic integrated circuits out of the research laboratory and into the marketplace has been made possible by the availability of high-quality epitaxial growth techniques for building up lattice-matched crystalline layers of indium gallium arsenide phosphide and indium phosphide (InGaAsP/InP). This III-V compound system is central to the light emitters and detectors used in the 1.3-micrometre and 1.5-micrometre wavelength ranges at which optical fibre has very low transmission loss.

Optical transmission. As the rates of transmission are increased from millions of bits (megabits) per second to billions of bits (gigabits) per second, commercially avail-

Molecular-beam epitaxy

Quantum wells

Band gap

able lasers encounter a physical limitation called "chirping," in which the optical frequency of the laser begins to waver during a pulse. Future systems, which may require from 2.4 to 30 gigabits per second, are probably going to be based on the use of a continuously operating distributed-feedback laser, whose output will be modulated in intensity by passing it through a modulator. This device consists of a crystal substrate of lithium niobate onto which a titanium channel is diffused to function as a light guide. The signal is encoded onto the light beam via a microwave radio-frequency feed through neighbouring channels in the coupler. Such a device is used only at the transmitter end of the optical path.

Both communications and computer systems rely on silica glass fibres to transmit light signals from lasers and LEDs. For long-distance transmission, optical-fibre cables are usually equipped with electro-optical repeater assemblies approximately every 100 kilometres. A new approach, called optical amplifiers, has been developed for deployment in transoceanic fibre-optic cables. Unlike traditional repeaters, optical amplifiers work by adding photons to a light signal without changing it to an electrical signal and without changing its bit-rate. Since they can be used at any desired transmission bit-rate, a transoceanic cable equipped with these devices can be upgraded to higher bit-rates simply by changing the lasers and photodiodes at each end. No retrofitting of higher bit-rate amplifiers is necessary.

Amplifying
optical
transmis-
sions

The optical amplifier is a module containing a semiconductor pump laser and a short length of optical fibre whose core has been doped with less than 0.1 percent erbium, an optically active rare-earth element. The pump laser is powered by an electrical conductor that runs the length of the cable. The amplifier functions by converting the optical energy generated by the pump source into signal photon energy. When a signal-carrying stream of laser pulses passes through the optical amplifier, it is combined with the pump light through a wavelength division multiplexer located in the module. The combined signal is fed through the erbium-doped fibre length, where the excited erbium ions contribute photons coherently to the signal. The amplified signal is then fed to the next section of cable for transmission to the next optical amplifier, perhaps 200 to 300 kilometres away.

(C.K.N.P.)

Materials for medicine

The treatment of many human disease conditions requires surgical intervention in order to assist, augment, sustain, or replace a diseased organ, and such procedures involve the use of materials foreign to the body. These materials, known as biomaterials, include synthetic polymers and, to a lesser extent, biological polymers, metals, and ceramics. Specific applications of biomaterials range from high-volume products such as blood bags, syringes, and needles to more challenging implantable devices designed to augment or replace a diseased human organ. The latter devices are used in cardiovascular, orthopedic, and dental applications as well as in a wide range of invasive treatment and diagnostic systems. Many of these devices have made possible notable clinical successes. For example, in cardiovascular applications, thousands of lives have been saved by heart valves, heart pacemakers, and large-diameter vascular grafts, and orthopedic hip-joint replacements have shown great long-term success in the treatment of patients suffering from debilitating joint diseases. With such a tremendous increase in medical applications, demand for a wide range of biomaterials grows by 5 to 15 percent each year. In the United States the annual market for surgical implants exceeds \$10 billion, approximately 10 percent of world demand.

Bio-
compati-
bility

Nevertheless, applications of biomaterials are limited by biocompatibility, the problem of adverse interactions arising at the junction between the biomaterial and the host tissue. Optimizing the interactions that occur at the surface of implanted biomaterials represents the most significant key to further advances, and an excellent basis for these advances can be found in the growing understanding of complex biological materials and in the development of

novel biomaterials custom-designed at the molecular level for specific medical applications.

This section describes biomaterials that are used in medicine, with emphasis on polymer materials and on the challenges associated with implantable devices used in the cardiovascular and orthopedic areas.

GENERAL REQUIREMENTS OF BIOMATERIALS

Research on developing new biomaterials is an interdisciplinary effort, often involving collaboration among materials scientists and engineers, biomedical engineers, pathologists, and clinicians to solve clinical problems. The design or selection of a specific biomaterial depends on the relative importance of the various properties that are required for the intended medical application. Physical properties that are generally considered include hardness, tensile strength, modulus, and elongation; fatigue strength, which is determined by a material's response to cyclic loads or strains; impact properties; resistance to abrasion and wear; long-term dimensional stability, which is described by a material's viscoelastic properties; swelling in aqueous media; and permeability to gases, water, and small biomolecules. In addition, biomaterials are exposed to human tissues and fluids, so that predicting the results of possible interactions between host and material is an important and unique consideration in using synthetic materials in medicine. Two particularly important issues in biocompatibility are thrombosis, which involves blood coagulation and the adhesion of blood platelets to biomaterial surfaces, and the fibrous-tissue encapsulation of biomaterials that are implanted in soft tissues.

Poor selection of materials can lead to clinical problems. One example of this situation was the choice of silicone rubber as a poppet in an early heart valve design. The silicone absorbed lipid from plasma and swelled sufficiently to become trapped between the metal struts of the valve. Another unfortunate choice as a biomaterial was Teflon (trademark), which is noted for its low coefficient of friction and its chemical inertness but which has relatively poor abrasion resistance. Thus, as an occluder in a heart valve or as an acetabular cup in a hip-joint prosthesis, Teflon may eventually wear to such an extent that the device would fail. In addition, degradable polyesterurethane foam was abandoned as a fixation patch for breast prostheses, because it offered a distinct possibility for the release of carcinogenic by-products as it degraded.

Besides their constituent polymer molecules, synthetic biomaterials may contain several additives, such as unreacted monomers and catalysts, inorganic fillers or organic plasticizers, antioxidants and stabilizers, and processing lubricants or mold-release agents on the material's surface. In addition, several degradation products may result from the processing, sterilization, storage, and ultimately implantation of a device. Many additives are beneficial—for example, the silica filler that is indispensable in silicone rubber for good mechanical performance or the antioxidants and stabilizers that prevent premature oxidative degradation of polyetherurethanes. Other additives, such as pigments, can be eliminated from biomedical products. Indeed, a "medical-grade" biomaterial is one that has had nonessential additives and potential contaminants excluded or eliminated from the polymer. In order to achieve this grade, the polymer may need to be solvent-extracted before use, thereby eliminating low-molecular-weight materials. Generally, additives in polymers are regarded with extreme suspicion, because it is often the additives rather than the constituent polymer molecules that are the source of adverse biocompatibility.

The
problem of
additives

POLYMER BIOMATERIALS

The majority of biomaterials used in humans are synthetic polymers such as the polyurethanes or Dacron (trademark: chemical name polyethylene terephthalate), rather than polymers of biological origin such as proteins or polysaccharides. The common synthetic biomaterials and their applications are listed in Table 2. Their properties vary widely, from the soft and delicate water-absorbing hydrogels made into contact lenses to the resilient elastomers found in short- and long-term cardiovascular devices or

Table 2: Biomaterials and Their Applications

material and application	biomedical device
Cardiovascular applications	
Polypropylene (isotactic)	prosthetic heart valve structures; sutures
Polytetrafluoroethylene	vascular grafts; catheter coating
Poly(2-hydroxyethyl methacrylate)	coatings for vascular grafts and catheters
Polyethylene terephthalate	vascular grafts, shunts
Polysulfones	prosthetic heart valves; artificial heart structures
Polyetherurethanes	intra-aortic balloons, catheters, percutaneous leads
Polyetherurethane ureas	artificial heart components; heart valves; device coatings
Soft-tissue applications	
Low-density polyethylene	tubings; shunts; syringes
Polytetrafluoroethylene	facial prostheses
Polymethyl methacrylate	intraocular lenses
Polyamides	sutures
Poly(2-hydroxyethyl methacrylate)	controlled drug release; contact lenses; artificial organs
Polyethylene terephthalate	tissue patches
Poly(lactide coglycolide)	bioresorbable sutures; controlled drug release
Polyorthoesters	controlled drug release
Polyanhydrides	controlled drug release
Polydimethylsiloxane	reconstructive devices; breast prostheses; shunts
Polyetherurethanes	wound dressing
Orthopedic/orthodontic applications	
High-density polyethylene	acetabular cups
Polymethyl methacrylate	dentures; bone cement; middle-ear prostheses
Extracorporeal applications	
Polypropylene (isotactic)	plasmapheresis membrane
Polytetrafluoroethylene	oxygenator membrane
Polyvinyl chloride	plasmapheresis membrane; blood bag
Polyacrylonitrile	hemodialysis membrane
Polydimethylsiloxane	oxygenator membrane

the high-strength acrylics used in orthopedics and dentistry. The properties of any material are governed by its chemical composition and by the intra- and intermolecular forces that dictate its molecular organization. Macromolecular structure in turn affects macroscopic properties and, ultimately, the interfacial behaviour of the material in contact with blood or host tissues.

Since the properties of each material are dependent on the chemical structure and macromolecular organization of its polymer chains, an understanding of some common structural features of various polymers provides considerable insight into their properties. Compared with complex biological molecules, synthetic polymers are relatively simple; often they comprise only one type of repeating subunit, analogous to a polypeptide consisting of just one repeating amino acid. On the basis of common structures and properties, synthetic polymers are classified into one of three categories: elastomers, which include natural and synthetic rubbers; thermoplastics; and thermosets. The properties that provide the basis for this classification include molecular weight, cross-link density, percent crystallinity, thermal transition temperature, and bulk mechanical properties.

Elastomers. Elastomers, which include rubber materials, have found wide use as biomaterials in cardiovascular and soft-tissue applications owing to their high elasticity, impact resistance, and gas permeability. Applications of elastomers include flexible tubing for pacemaker leads, vascular grafts, and catheters; biocompatible coatings and pumping diaphragms for artificial hearts and left-ventricular assist devices; grafts for reconstructive surgery and maxillofacial operations; wound dressings; breast prostheses; and membranes for implantable biosensors.

Elastomers are typically amorphous with low cross-link density (although linear polyurethane block copolymers are an important exception). This gives them low to moderate modulus and tensile properties as well as high elasticity. For example, elastomeric devices can be extended by 100 to 1,000 percent of their initial dimensions without causing any permanent deformation to the material. Silicone rubbers such as Silastic (trademark), produced by the American manufacturer Dow Corning, Inc., are cross-linked, so that they cannot be melted or dissolved—although swelling may occur in the presence of a good solvent. Such properties contrast with those of the linear polyurethane elastomers, which consist of soft polyether amorphous segments and hard urethane-containing glassy or crystalline segments. The two segments are incompatible at room temperature and undergo microphase separation, forming hard domains dispersed in an amorphous

matrix. A key feature of this macromolecular organization is that the hard domains serve as physical cross-links and reinforcing filler. This results in elastomeric materials that possess relatively high modulus and extraordinary long-term stability under sustained cyclic loading. In addition, they can be processed by methods common to thermoplastics.

Thermoplastics. Many common thermoplastics, such as polyethylene and polyester, are used as biomaterials. Thermoplastics usually exhibit moderate to high tensile strength (5 to 1,000 megapascals) with moderate elongation (2 to 100 percent), and they undergo plastic deformation at high strains. Thermoplastics consist of linear or branched polymer chains; consequently, most can undergo reversible melt-solid transformation on heating, which allows for relatively easy processing or reprocessing. Depending on the structure and molecular organization of the polymer chains, thermoplastics may be amorphous (*e.g.*, polystyrene), semicrystalline (*e.g.*, low-density polyethylene), or highly crystalline (*e.g.*, high-density polyethylene), or they may be processed into highly crystalline textile fibres (*e.g.*, polyester Dacron).

Some thermoplastic biomaterials, such as polylactic acid and polyglycolic acid, are polymers based on a repeating amino acid subunit. These polypeptides are biodegradable, and, along with biodegradable polyesters and polyorthoesters, they have applications in absorbable sutures and drug-release systems. The rate of biodegradation in the body can be adjusted by using copolymers. These are polymers that link two different monomer subunits into a single polymer chain. The resultant biomaterial exhibits properties, including biodegradation, that are intermediate between the two homopolymers.

Thermosets. Thermosetting polymers find only limited application in medicine, but their characteristic properties, which combine high strength and chemical resistance, are useful for some orthopedic and dental devices. Thermosetting polymers such as epoxies and acrylics are chemically inert, and they also have high modulus and tensile properties with negligible elongation (1 to 2 percent). The polymer chains in these materials are highly cross-linked and therefore have severely restricted macromolecular mobility; this limits extension of the polymer chains under an applied load. As a result, thermosets are strong but brittle materials.

Cross-linking inhibits close packing of polymer chains, preventing formation of crystalline regions. Another consequence of extensive cross-linking is that thermosets do not undergo solid-melt transformation on heating, so that they cannot be melted or reprocessed.

Elasticity,
impact
resistance,
and gas
perme-
ability

Strength
and
chemical
resistance

APPLICATIONS OF BIOMATERIALS

Cardiovascular devices. Biomaterials are used in many blood-contacting devices. These include artificial heart valves, synthetic vascular grafts, ventricular assist devices, drug-release systems, extracorporeal systems, and a wide range of invasive treatment and diagnostic systems. An important issue in the design and selection of materials is the hemodynamic conditions in the vicinity of the device. For example, mechanical heart valve implants are intended for long-term use. Consequently, the hinge points of each valve leaflet and the materials must have excellent wear and fatigue resistance in order to open and close 80 times per minute for many years after implantation. In addition, the open valve must minimize disturbances to blood flow as blood passes from the left ventricle of the heart, through the heart valve, and into the ascending aorta of the arterial vascular system. To this end, the bileaflet valve disks of one type of implant are coated with pyrolytic carbon, which provides a relatively smooth, chemically inert surface. This is an important property, because surface roughness will cause turbulence in the blood flow, which in turn may lead to hemolysis of red cells, provide sites for adventitious bacterial adhesion and subsequent colonization, and, in areas of blood stasis, promote thrombosis and blood coagulation. The carbon-coated holding ring of this implant is covered with Dacron mesh fabric so that the surgeon can sew and fix the device to adjacent cardiac tissues. Furthermore, the porous structure of the Dacron mesh promotes tissue integration, which occurs over a period of weeks after implantation.

While the possibility of thrombosis can be minimized in blood-contacting biomaterials, it cannot be eliminated entirely. For this reason, patients who receive artificial heart valves or other blood-contacting devices also receive anticoagulation therapy. This is needed because all foreign surfaces initiate blood coagulation and platelet adhesion to some extent. Platelets are circulating cellular components of blood, two to four micrometres in size, that attach to foreign surfaces and actively participate in blood coagulation and thrombus formation. Research on new biomaterials for cardiovascular applications is largely devoted to understanding thrombus formation and to developing novel surfaces for biomaterials that will provide improved blood compatibility.

Synthetic vascular graft materials are used to patch injured or diseased areas of arteries, for replacement of whole segments of larger arteries such as the aorta, and for use as sewing cuffs (as with the heart valve mentioned above). Such materials need to be flexible to allow for the difficulties of implantation and to avoid irritating adjacent tissues; also, the internal diameter of the graft should remain constant under a wide range of flexing and bending conditions, and the modulus or compliance of the vessel should be similar to that of the natural vessel. These aims are largely achieved by crimped woven Dacron and expanded polytetrafluoroethylene (ePTFE). Crimping of Dacron in processing results in a porous vascular graft that may be bent 180° or twisted without collapsing the internal diameter.

A biomaterial used for blood vessel replacement will be in contact not only with blood but also with adjacent soft tissues. Experience with different materials has shown that tissue growth into the interstices of the biomaterials aids healing and integration of the material with host tissue after implantation. In order for the tissue, which consists mostly of collagen, to grow in the graft, the vascular graft must have an open structure with pores at least 10 micrometres in diameter. These pores allow new blood capillaries that develop during healing to grow into the graft, and the blood then provides oxygen and other nutrients for fibroblasts and other cells to survive in the biomaterial matrix. Fibroblasts synthesize the structural protein tropocollagen, which is needed in the development of new fibrous tissue as part of the healing response to a surgical wound.

Occasionally, excessive tissue growth may be observed at the anastomosis, which is where the graft is sewn to the native artery. This is referred to as internal hyperplasia and is thought to result from differences in compliance

between the graft and the host vessels. In addition, in order to optimize compatibility of the biomaterial with the blood, the synthetic graft eventually should be coated with a confluent layer of host endothelial cells, but this does not occur with current materials. Therefore, most proposed modifications to existing graft materials involve potential improvements in blood compatibility.

Artificial heart valves and vascular grafts, while not ideal, have been used successfully and have saved many thousands of lives. However, the risk of thrombosis has limited the success of existing cardiovascular devices and has restricted potential application of the biomaterials to other devices. For example, there is an urgent clinical need for blood-compatible, synthetic vascular grafts of small diameter in peripheral vascular surgery—*e.g.*, in the legs—but this is currently impracticable with existing biomaterials because of the high risk of thrombotic occlusion. Similarly, progress with implantable miniature sensors, designed to measure a wide range of blood conditions continuously, has been impeded because of problems directly attributable to the failure of existing biomaterials. With such biocompatibility problems resolved, biomedical sensors would provide a very important contribution to medical diagnosis and monitoring. Considerable advances have been made in the ability to manipulate molecular architecture at the surfaces of materials by using chemisorbed or physisorbed monolayer films. Such progress in surface modification, combined with the development of nanoscale probes that permit examination at the molecular and submolecular level, provide a strong basis for optimism in the development of specialty biomaterials with improved blood compatibility.

Orthopedic devices. Joint replacements, particularly at the hip, and bone fixation devices have become very successful applications of materials in medicine. The use of pins, plates, and screws for bone fixation to aid recovery of bone fractures has become routine, with the number of annual procedures approaching five million in the United States alone. In joint replacement, typical patients are age 55 or older and suffer from debilitating rheumatoid arthritis, osteoarthritis, or osteoporosis. Orthopedic surgeries for artificial joints exceed 1.5 million each year, with actual joint replacement accounting for about half of the procedures. A major focus of research is the development of new biomaterials for artificial joints intended for younger, more active patients.

Hip-joint replacements are principally used for structural support. Consequently, they are dominated by materials that possess high strength, such as metals, tough plastics, and reinforced polymer-matrix composites. In addition, biomaterials used for orthopedic applications must have high modulus, long-term dimensional stability, high fatigue resistance, long-term biostability, excellent abrasion resistance, and biocompatibility (*i.e.*, there should be no adverse tissue response to the implanted device). Early developments in this field used readily available materials such as stainless steels, but evidence of corrosion after implantation led to their replacement by more stable materials, particularly titanium alloys, cobalt-chromium-molybdenum alloys, and carbon fibre-reinforced polymer composites. A typical modern artificial hip consists of a nitrided and highly polished cobalt-chromium ball connected to a titanium alloy stem that is inserted into the femur and cemented into place by in situ polymerization of polymethylmethacrylate. The articulating component of the joint consists of an acetabular cup made of tough, creep-resistant, ultrahigh-molecular-weight polyethylene. Abrasion at the ball-and-cup interface can lead to the production of wear particles, which in turn can lead to significant inflammatory reaction by the host. Consequently, much research on the development of hip-joint materials has been devoted to optimizing the properties of the articulating components in order to eliminate surface wear. Other modifications include porous coatings made by sintering the metal surface or coatings of wire mesh or hydroxyapatite; these promote bone growth and integration between the implant and the host, eliminating the need for an acrylic bone cement.

While the strength of the biomaterials is important, an-

The need for improved blood compatibility

Vascular graft materials

The artificial hip

other goal is to match the mechanical properties of the implant materials with those of the bone in order to provide a uniform distribution of stresses (load sharing). If a bone is loaded insufficiently, the stress distribution will be made asymmetric, and this will lead to adaptive remodeling with cortical thinning and increased porosity of the bone. Such lessons in structure hierarchy and in the structure-property relationships of materials have been obtained from studies on biologic composite materials, and they are being translated into new classes of synthetic biomaterials. One development is carbon fibre-reinforced polymer-matrix composites. Typical matrix polymers include polysulfone and polyetheretherketones. The strength of these composites is lower than that of metals, but it more closely approximates that of bone. (R.E.M.)

BIBLIOGRAPHY

General works. Overviews of the properties and production of all engineering materials can be found in the following texts: JAMES F. SHACKELFORD, *Introduction to Materials Science for Engineers*, 3rd ed. (1992); WILLIAM D. CALLISTER, JR., *Materials Science and Engineering: An Introduction*, 2nd ed. (1991); RICHARD A. FLINN and PAUL K. TROJAN, *Engineering Materials and Their Applications*, 4th ed. (1990); DONALD R. ASKELAND, *The Science and Engineering of Materials*, 2nd ed. (1989); and MICHAEL F. ASHBY and DAVID R.H. JONES, *Engineering Materials: An Introduction to Their Properties and Applications* (1980), readable even for those with no previous materials science background and providing clearly described examples of innovative ways to use materials. *Materials Science and Engineering for the 1990s* (1989) comprehensively describes new directions to be taken in materials science; it is written in a readily comprehensible manner by members of committees of the National Research Council (U.S.). An entire issue of *Advanced Materials & Processes*, vol. 141, no. 1 (January 1992), is devoted to a forecast of developments in various materials, trends in materials processing, and advances in testing and characterization of materials. MICHAEL B. BEVER (ed.), *Encyclopedia of Materials Science and Engineering*, 8 vol. (1986), with supplementary vols., is a comprehensive reference work.

Materials for energy. Three journal articles on the topic are RICHARD S. CLAASEN and LOUIS A. GIRIFALCO, "Materials for Energy Utilization," *Scientific American*, 255(4):102-104, 109-112, 117 (October 1986); RICHARD S. CLAASEN, "Materials for Advanced Energy Technologies," *Science*, 191(4227):739-745 (Feb. 20, 1976); and BERNARD L. COHEN, "The Disposal of Radioactive Wastes from Fission Reactors," *Scientific American*, 236(6):21-31 (June 1977).

Materials for ground transportation. A lucid account of the shift away from conventional steels in modern automobiles is found in the excellent introductory article by W. DALE COMPTON and NORMAN A. GJOSTEIN, "Materials for Ground Transportation," *Scientific American*, 255(4):92-100 (October 1986). KAREN WRIGHT, "The Shape of Things to Go," *Scientific American*, 262(5):92-101 (May 1990), projects the effect of advanced technology on automobiles of the future.

Materials for aerospace. An overview is found in MORRIS A. STEINBERG, "Materials for Aerospace," *Scientific American*, 255(4):66-72 (October 1986). An entire issue of *Advanced Materials & Processes*, vol. 137, no. 4 (April 1990), is devoted to aerospace materials and applications. A special section on frontiers in materials science in *Science*, 255(5048):1077-1112 (Feb. 28, 1992), discusses polymers and aircraft engine materials, among other subjects.

Works on composites include TSU-WEI CHOU, ROY L. MCCULLOUGH, and R. BYRON PIPES, "Composites," *Scientific American*, 255(4):192-203 (October 1986); ROY L. MCCULLOUGH, *Concepts*

of Fiber-Resin Composites (1971); STEPHEN W. TSAI and H. THOMAS HAHN, *Introduction to Composite Materials* (1980); and JACK R. VINSON and TSU-WEI CHOU, *Composite Materials and Their Use in Structures* (1975).

Materials for communications. A useful introduction is by JOHN S. MAYO, "Materials for Information and Communication," *Scientific American*, 255(4):58-66 (October 1986). Discussions of electronic and photonic materials may be found in the following essays, all from *AT&T Technical Journal*: in vol. 69, no. 6 (November/December 1990), see C. KUMAR N. PATEL, "Materials and Processing: Core Competencies and Strategic Resources," pp. 2-8; KENNETH E. BENSON, LIONEL C. KIMERLING, and PETER T. PANOUSIS, "Reaching the Limits in Silicon Processing," pp. 16-31; ELSA REICHMANIS and LARRY F. THOMPSON, "Challenges in Lithographic Materials and Processes," pp. 32-45; and JAMES W. MITCHELL, JORGE LUIS VALDES, and GARDY CADET, "Benign Precursors for Semiconductor Processing," pp. 101-112; in vol. 68, no. 1 (January/February 1989), see JIM E. CLEMANS *et al.*, "Bulk III-V Compound Semiconductor Crystal Growth," pp. 29-42; and W. DEXTER JOHNSTON, JR., MICHAEL A. DIGUISEPPE, and DANIEL P. WILT, "Liquid and Vapor Phase Growth of III-V Materials for Photonic Devices," pp. 53-63; and in vol. 71, no. 1 (January/February 1992), see JOHN L. ZYSKIND *et al.*, "Erbium-Doped Fiber Amplifiers and the Next Generation of Lightwave Systems," pp. 53-62.

Materials for medicine. ROBERT A. FULLER and JONATHAN J. ROSEN, "Materials for Medicine," *Scientific American*, 255(4):118-125 (October 1986), offers an overview of the subject. S.A. BARENBERG, "Abridged Report of the Committee to Survey the Needs and Opportunities for the Biomaterials Industry," *Journal of Biomedical Materials Research*, 22(12):1267-92 (December 1988), surveys the applications of materials in medicine and highlights projected areas of clinical need. JOON BU PARK, *Biomaterials Science and Engineering* (1984), provides a qualitative university-level introduction to the field of biomaterials. MICHAEL SZYCHER (ed.), *Biocompatible Polymers, Metals, and Composites* (1983), is a collection of detailed review articles that covers materials in medicine and biocompatibility and contains a pragmatic assessment of clinical and commercial aspects, including how to sterilize and package biomaterials. HARRY R. ALLCOCK and FREDERICK W. LAMPE, *Contemporary Polymer Chemistry*, 2nd ed. (1990), is a basic textbook of polymer science, providing a university-level introduction to synthesis and characterization of polymers, including biomedical polymers. Advanced biomaterials texts with emphasis on research include MICHAEL SZYCHER (ed.), *High Performance Biomaterials: A Comprehensive Guide to Medical and Pharmaceutical Applications* (1991), research articles covering orthopedic and cardiovascular biomaterials as well as most other areas of materials in medicine; JOSEPH D. ANDRADE (ed.), *Surface and Interfacial Aspects of Biomedical Polymers*, vol. 1, *Surface Chemistry and Physics* (1985), articles on surface characterization methods applied to biomaterials, including a quantitative presentation of the interactions of blood components (especially proteins) with biomaterial surfaces; HOWARD P. GREISLER, *New Biologic and Synthetic Vascular Prostheses* (1991), a biological perspective on blood interactions, wound healing, and tissue integration with biomaterials and surface modified materials; D.F. WILLIAMS, *Blood Compatibility*, 2 vol. (1987), detailed review articles covering blood interactions with biomaterials and prosthetic devices and methods of modifying the surface of biomaterials; and I.L. GOLDSMITH and V.T. TURITTO, "Rheological Aspects of Thrombosis and Hemostasis: Basic Principles and Applications," *Thrombosis and Haemostasis*, 55(3):415-435 (1986), a detailed and quantitative review article that describes and models blood flow and rheology in the vascular system, including the effects of different blood components.

(L.A.G./J.D.V./R.L.McC./D.S.K./C.K.N.P./R.E.M.)

The Foundations of Mathematics

The study of the foundations of mathematics is the examination of the underlying assumptions and procedures of mathematics, and its limitations. Because mathematics has served as a model for rational inquiry in the West and is used extensively in the sciences, foundational studies have far-reaching consequences for the reliability and extensibility of rational thought itself.

For 2,000 years the foundations of mathematics seemed perfectly solid. Euclid's *Elements* (c. 300 BC), which presented a set of formal logical arguments based on a few basic terms and axioms, provided a systematic method of rational exploration that guided mathematicians, philosophers, and scientists well into the 19th century. Even serious objections to the lack of rigour in Sir Isaac Newton's notion of fluxions (derivatives) in the calculus, raised by the Anglo-Irish empiricist George Berkeley (among others), did not call into question the basic foundations of mathematics. The discovery in the 19th century of consistent alternative geometries, however, precipitated a crisis, for it showed that Euclidean geometry, based on seemingly the most intuitively obvious axiomatic assumptions, did not correspond with reality as mathematicians had believed.

This, together with the bold discoveries of the German mathematician Georg Cantor in set theory, made it clear that, to avoid further confusion and satisfactorily answer paradoxical results, a new and more rigorous foundation for mathematics was necessary.

Thus began the 20th-century quest to rebuild mathematics on a new basis independent of geometric intuitions. Early efforts included those of the logicist school of the British mathematicians Bertrand Russell and Alfred North Whitehead, the formalist school of the German mathematician David Hilbert, the intuitionist school of the Dutch mathematician L.E.J. Brouwer, and the French set theory school of mathematicians collectively writing under the pseudonym of Nicolas Bourbaki. Some of the most promising current research is based on the development of category theory by the American mathematician Saunders Mac Lane and the Polish-born American mathematician Samuel Eilenberg following World War II.

This article presents the historical background of foundational questions and 20th-century efforts to construct a new foundational basis for mathematics.

This article is divided into the following sections:

Ancient Greece to the Enlightenment 566

- Arithmetic or geometry
- Being versus becoming
- Universals
- The axiomatic method
- Number systems

The reexamination of infinity 568

- Calculus reopens foundational questions
- Non-Euclidean geometries
- Cantor

The quest for rigour 569

- Formal foundations 569
 - Set theoretic beginnings
 - Foundational logic
 - Impredicative constructions
 - Nonconstructive arguments
 - Intuitionistic logic

Other logics

- Formalism
- Gödel
- Recursive definitions
- Computers and proof

Category theory 572

- Abstraction in mathematics
- Isomorphic structures
- Topos theory
- Intuitionistic type theories
- Internal language
- Gödel and category theory
- The search for a distinguished model
- Boolean local topoi
- One distinguished model or many models

Bibliography 574

Ancient Greece to the Enlightenment

A remarkable amount of practical mathematics, some of it even fairly sophisticated, was already developed as early as 2000 BC by the agricultural civilizations of Egypt and Mesopotamia, and perhaps even farther east. However, the first to exhibit an interest in the foundations of mathematics were the ancient Greeks.

Arithmetic or geometry. Early Greek philosophy was dominated by a dispute as to which is more basic, arithmetic or geometry, and thus whether mathematics should be concerned primarily with the (positive) integers or the (positive) reals, the latter then being conceived as ratios of geometric quantities. (The Greeks confined themselves to positive numbers, as negative numbers were introduced only much later in India by Brahmagupta.) Underlying this dispute was a perceived basic dichotomy, not confined to mathematics but pervading all nature: is the universe made up of discrete atoms (as the philosopher Democritus believed) which hence can be counted, or does it consist of one or more continuous substances (as Thales of Miletus is reputed to have believed) and thus can only be measured? This dichotomy was presumably inspired by a linguistic distinction, analogous to that between English count nouns, such as "apple," and mass nouns, such as "water." As Aristotle later pointed out, in an effort to mediate between these divergent positions, water can be measured by counting cups.

The Pythagorean school of mathematics, founded on the

doctrines of the Greek philosopher Pythagoras, originally insisted that only natural and rational numbers exist. Its members only reluctantly accepted the discovery that $\sqrt{2}$, the ratio of the diagonal of a square to its side, could not be expressed as the ratio of whole numbers. The remarkable proof of this fact has been preserved by Aristotle. (See ANALYSIS: *Real analysis*.)

The contradiction between rationals and reals was finally resolved by Eudoxus of Cnidus, a disciple of Plato, who pointed out that two ratios of geometric quantities are equal if and only if they partition the set of (positive) rationals in the same way, thus anticipating the German mathematician Richard Dedekind (1831–1916), who defined real numbers as such partitions.

Being versus becoming. Another dispute among pre-Socratic philosophers was more concerned with the physical world. Parmenides claimed that in the real world there is no such thing as change and that the flow of time is an illusion, a view with parallels in the Einstein-Minkowski four-dimensional space-time model of the universe. Heraclitus, on the other hand, asserted that change is all-pervasive and is reputed to have said that one cannot step into the same river twice.

Zeno of Elea, a follower of Parmenides, claimed that change is actually impossible and produced four paradoxes to show this. The most famous of these describes a race between Achilles and a tortoise. Since Achilles can run much faster than the tortoise, let us say twice as fast, the latter is allowed a head start of one mile. When Achilles

Discrete
atoms
versus
continuous
substances

has run one mile, the tortoise will have run half as far again—that is, half a mile. When Achilles has covered that additional half-mile, the tortoise will have run a further quarter-mile. After $n + 1$ stages, Achilles has run

$$1 + \frac{1}{2} + \dots + \frac{1}{2^n} = 2 - \frac{1}{2^n}$$

miles and the tortoise has run

$$1 + \frac{1}{2} + \dots + \frac{1}{2^n} + \frac{1}{2^{n+1}}$$

miles, being still $1/2^{n+1}$ miles ahead. So how can Achilles ever catch up with the tortoise (Figure 1)?

Zeno's paradoxes may also be interpreted as showing that space and time are not made up of discrete atoms but are substances which are infinitely divisible. Mathematically speaking, his argument involves the sum of the infinite geometric progression

$$1 + \frac{1}{2} + \frac{1}{4} + \dots,$$

no finite partial sum of which adds up to 2. As Aristotle would later say, this progression is only potentially infinite. It is now understood that Zeno was trying to come to grips with the notion of limit, which was not formally explained until the 19th century, although a start in that direction had been made by the French encyclopaedist Jean Le Rond d'Alembert (1717–83).

Universals. The Athenian philosopher Plato believed that mathematical entities are not just human inventions but have a real existence. For instance, according to Plato, the number 2 is an ideal object. This is sometimes called an "idea," from the Greek *eide*, or "universal," from the Latin *universalis*, meaning "that which pertains to all." But Plato did not have in mind a "mental image," as "idea" is usually used. The number 2 is to be distinguished from a collection of two stones or two apples or, for that matter, two platinum balls in Paris.

What, then, are these Platonic ideas? Already in ancient Alexandria some people speculated that they are words. This is why the Greek word *logos*, originally meaning "word," later acquired a theological meaning as denoting the ultimate reality behind the "thing." An intense debate occurred in the Middle Ages over the ontological status of universals. Three dominant views prevailed: realism, from the Latin *res* ("thing"), which asserts that universals have an extra-mental reality—that is, they exist independently of perception; conceptualism, which asserts that universals exist as entities within the mind but have no extra-mental existence; and nominalism, from the Latin *nomen* ("name"), which asserts that universals exist neither in the mind nor in the extra-mental realm but are merely names that refer to collections of individual objects.

It would seem that Plato believed in a notion of truth independent of the human mind. In the *Meno* Plato's teacher Socrates asserts that it is possible to come to know this truth by a process akin to memory retrieval. Thus, by clever questioning, Socrates managed to bring an uneducated person to "remember," or rather to reconstruct, the proof of a mathematical theorem.

The axiomatic method. Perhaps the most important contribution to the foundations of mathematics made by the ancient Greeks was the axiomatic method and the notion of proof. This was insisted upon in Plato's Academy and reached its high point in Alexandria about 300 BC with Euclid's *Elements*. This notion survives today, except for some cosmetic changes.

The idea is this: there are a number of basic mathematical truths, called axioms or postulates, from which other true statements may be derived in a finite number of steps. It may take considerable ingenuity to discover a proof; but it is now held that it must be possible to check mechanically, step by step, whether a purported proof is indeed correct, and nowadays a computer should be able to do this. The mathematical statements that can be proved are called theorems, and it follows that, in principle, a mechanical device, such as a modern computer, can generate all theorems.

Two questions about the axiomatic method were left

Infinite geometric progression

Euclid's *Elements*

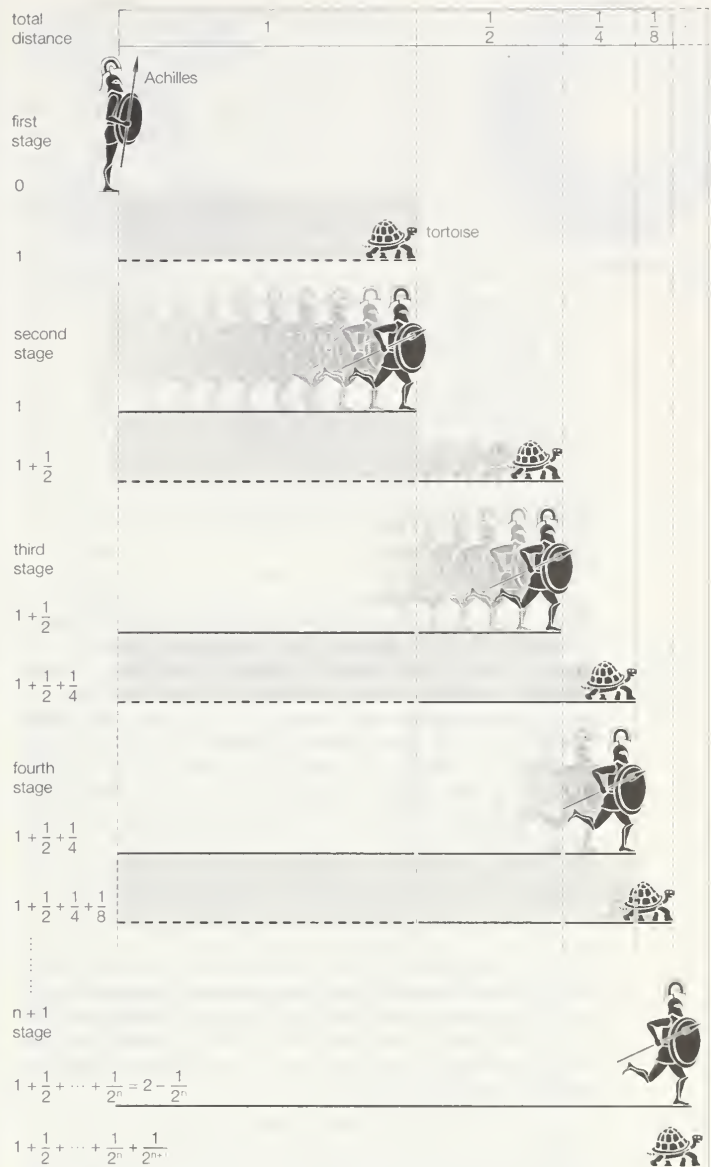


Figure 1: Zeno's paradox, illustrated by Achilles racing a tortoise.

Encyclopædia Britannica, Inc.

unanswered by the ancients: are all mathematical truths axioms or theorems (this is referred to as completeness), and can it be determined mechanically whether a given statement is a theorem (this is called decidability)? These questions were raised implicitly by David Hilbert (1862–1943) about 1900 and were resolved later in the negative, completeness by the Austrian-American logician Kurt Gödel (1906–78) and decidability by the American logician Alonzo Church (1903–95).

Euclid's work dealt with number theory and geometry, essentially all the mathematics then known. Since the middle of the 20th century a gradually changing group of mostly French mathematicians under the pseudonym Nicolas Bourbaki has tried to emulate Euclid in writing a new *Elements of Mathematics* based on their theory of structures. Unfortunately, they just missed out on the new ideas from category theory.

Number systems. While the ancient Greeks were familiar with the positive integers, rationals, and reals, zero (used as an actual number instead of denoting a missing number) and the negative numbers were first used in India, as far as is known, by Brahmagupta in the 7th century AD. Complex numbers were introduced by the Italian Renaissance mathematician and physician Gerolamo Cardano (1501–76), not just to solve equations such as $x^2 + 1 = 0$ but because they were needed to find real solutions of certain cubic equations with real coefficients.

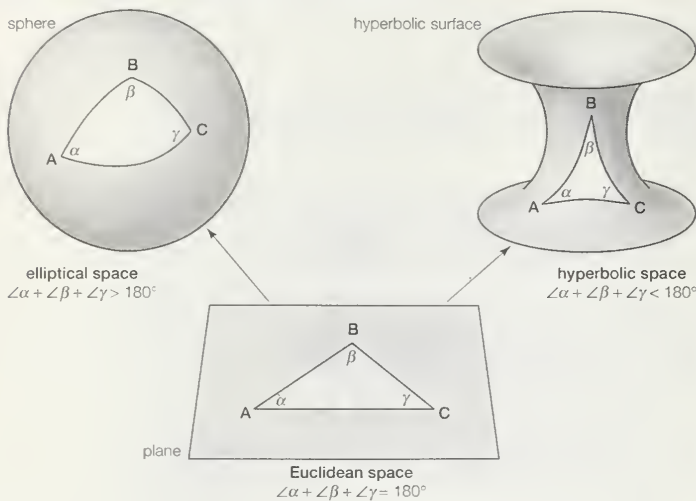


Figure 2: Contrasting triangles in Euclidean, elliptic, and hyperbolic spaces.

Encyclopædia Britannica, Inc

Much later, the German mathematician Carl Friedrich Gauss (1777–1855) proved the fundamental theorem of algebra, that all equations with complex coefficients have complex solutions, thus removing the principal motivation for introducing new numbers. Still, the Irish mathematician Sir William Rowan Hamilton (1805–65) and the French mathematician Olinde Rodrigues (1794–1851) invented quaternions in the mid-19th century, but these proved to be less popular in the scientific community until quite recently.

Currently, a logical presentation of the number system, as taught at the university level, would be as follows:

$$\mathbf{N} \rightarrow \mathbf{Z} \rightarrow \mathbf{Q} \rightarrow \mathbf{R} \rightarrow \mathbf{C} \rightarrow \mathbf{H}.$$

Here the boldfaced letters, introduced by Nicolas Bourbaki, refer to the natural numbers, integers, rationals, reals, complex numbers, and quaternions, respectively, and the arrows indicate inclusion of each number system into the next. However, as has been shown, the historical development proceeds differently:

$$\mathbf{N}^+ \rightarrow \mathbf{Q}^+ \rightarrow \mathbf{R}^+ \rightarrow \mathbf{R} \rightarrow \mathbf{C} \rightarrow \mathbf{H},$$

where the plus sign indicates restriction to positive elements. This is the development, up to **R**, which is often adhered to at the high-school level.

The reexamination of infinity

Calculus reopens foundational questions. Although mathematics flourished after the end of the Classical Greek period for 800 years in Alexandria and, after an interlude in India and the Islāmic world, again in Renaissance Europe, philosophical questions concerning the foundations of mathematics were not raised until the invention of calculus and then not by mathematicians but by the philosopher George Berkeley (1685–1753).

Sir Isaac Newton in England and Gottfried Wilhelm Leibniz in Germany had independently developed the calculus on a basis of heuristic rules and methods markedly deficient in logical justification. As is the case in many new developments, utility outweighed rigour, and, though Newton's fluxions (or derivatives) and Leibniz's infinitesimals (or differentials) lacked a coherent rational explanation, their power in answering heretofore unanswerable questions was undeniable. Unlike Newton, who made little effort to explain and justify fluxions, Leibniz, as an eminent and highly regarded philosopher, was influential in propagating the idea of infinitesimals, which he described as infinitely small actual numbers—that is, less than $1/n$ in absolute value for each positive integer n and yet not equal to zero. Berkeley, concerned over the deterministic and atheistic implications of philosophical mechanism, set out to reveal contradictions in the calculus in his influential book *The Analyst; or, A Discourse Addressed to an*

Infidel Mathematician. There he scathingly wrote about these fluxions and infinitesimals, "They are neither finite quantities, nor quantities infinitely small, nor yet nothing. May we not call them the ghosts of departed quantities?" and further asked, "Whether mathematicians, who are so delicate in religious points, are strictly scrupulous in their own science? Whether they do not submit to authority, take things upon trust, and believe points inconceivable?"

Berkeley's criticism was not fully met until the 19th century, when it was realized that, in the expression dy/dx , dx and dy need not lead an independent existence. Rather, this expression could be defined as the limit of ordinary ratios $\Delta y/\Delta x$, as Δx approaches zero without ever being zero. Moreover, the notion of limit was then explained quite rigorously, in answer to such thinkers as Zeno and Berkeley. (See ANALYSIS: *Real analysis*.)

It was not until the middle of the 20th century that the logician Abraham Robinson (1918–74) showed that the notion of infinitesimal was in fact logically consistent and that, therefore, infinitesimals could be introduced as new kinds of numbers. This led to a novel way of presenting the calculus, called nonstandard analysis, which has, however, not become as widespread and influential as it might have.

Robinson's argument was this: if the assumptions behind the existence of an infinitesimal ξ led to a contradiction, then this contradiction must already be obtainable from a finite set of these assumptions, say from:

$$0 < \xi, \xi < 1, \xi < \frac{1}{2}, \dots, \xi < \frac{1}{n}.$$

But this finite set is consistent, as is seen by taking $\xi = 1/(n + 1)$.

Non-Euclidean geometries. When Euclid presented his axiomatic treatment of geometry, one of his assumptions, his fifth postulate, appeared to be less obvious or fundamental than the others. As it is now conventionally formulated, it asserts that there is exactly one parallel to a given line through a given point. Attempts to derive this from Euclid's other axioms did not succeed, and, at the beginning of the 19th century, it was realized that Euclid's fifth postulate is, in fact, independent of the others. It was then seen that Euclid had described not the one true geometry but only one of a number of possible geometries.

Elliptic and hyperbolic geometries. Within the framework of Euclid's other four postulates (and a few that he omitted), there were also possible elliptic and hyperbolic geometries. In plane elliptic geometry there are no parallels to a given line through a given point; it may be viewed as the geometry of a spherical surface on which antipodal points have been identified and all lines are great circles. This was not viewed as revolutionary. More exciting was plane hyperbolic geometry, developed independently by the Hungarian mathematician János Bolyai (1802–60) and the Russian mathematician Nikolay Lobachevsky (1792–1856), in which there is more than one parallel to a given line through a given point. This geometry is more difficult to visualize, but a helpful model presents the hyperbolic plane as the interior of a circle, in which straight lines take the form of arcs of circles perpendicular to the circumference.

Another way to distinguish the three geometries is to look at the sum of the angles of a triangle. It is 180° in Euclidean geometry, as first reputedly discovered by Thales of Miletus (fl. 6th century BC), whereas it is more than 180° in elliptic and less than 180° in hyperbolic geometry. See Figure 2.

Riemannian geometry. The discovery that there is more than one geometry was of foundational significance and contradicted the German philosopher Immanuel Kant (1724–1804). Kant had argued that there is only one true geometry, Euclidean, which is known to be true a priori by an inner faculty (or intuition) of the mind. For Kant, and practically all other philosophers and mathematicians of his time, this belief in the unassailable truth of Euclidean geometry formed the foundation and justification for further explorations into the nature of reality. With the discovery of consistent non-Euclidean geometries, there was a subsequent loss of certainty and trust in this innate

Non-standard analysis

Euclid's parallel postulate

Infinitesimals

intuition, and this was fundamental in separating mathematics from a rigid adherence to an external sensory order (no longer vouchsafed as “true”) and led to the growing abstraction of mathematics as a self-contained universe. This divorce from geometric intuition added impetus to later efforts to rebuild assurance of truth on the basis of logic.

What then is the correct geometry for describing the space (actually space-time) we live in? It turns out to be none of the above, but a more general kind of geometry, as was first discovered by the German mathematician Bernhard Riemann (1826–66). In the early 20th century, Albert Einstein showed, in the context of his general theory of relativity, that the true geometry of space is only approximately Euclidean. It is a form of Riemannian geometry in which space and time are linked in a four-dimensional manifold, and it is the curvature at each point that is responsible for the gravitational “force” at that point. Einstein spent the last part of his life trying to extend this idea to the electromagnetic force, hoping to reduce all physics to geometry, but a successful unified field theory eluded him.

Cantor. In the 19th century, the German mathematician Georg Cantor (1845–1918) returned once more to the notion of infinity and showed that, surprisingly, there is not just one kind of infinity but many kinds. In particular, while the set \mathbf{N} of natural numbers and the set of all subsets of \mathbf{N} are both infinite, the latter collection is more numerous, in a way that Cantor made precise, than the former. He proved that \mathbf{N} , \mathbf{Z} , and \mathbf{Q} all have the same size, since it is possible to put them into one-to-one correspondence with one another, but that \mathbf{R} is bigger, having the same size as the set of all subsets of \mathbf{N} .

However, Cantor was unable to prove the so-called continuum hypothesis, which asserts that there is no set that is larger than \mathbf{N} yet smaller than the set of its subsets. It was shown only in the 20th century, by Gödel and the American logician Paul Cohen (b. 1934), that the continuum hypothesis can be neither proved nor disproved from the usual axioms of set theory. Cantor had his detractors, most notably the German mathematician Leopold Kronecker (1823–91), who felt that Cantor’s theory was too metaphysical and that his methods were not sufficiently constructive (see below *The quest for rigour: Formal foundations: Nonconstructive arguments*).

The quest for rigour

FORMAL FOUNDATIONS

Set theoretic beginnings. While laying rigorous foundations for mathematics, 19th-century mathematicians discovered that the language of mathematics could be reduced to that of set theory (developed by Cantor), dealing with membership (\in) and equality ($=$), together with some rudimentary arithmetic, containing at least symbols for zero (0) and successor (S). Underlying all this were the basic logical concepts: conjunction (\wedge), disjunction (\vee), implication (\supset), negation (\neg), and the universal (\forall) and existential (\exists) quantifiers (formalized by the German mathematician Gottlob Frege [1848–1925]). (The modern notation owes more to the influence of the English logician Bertrand Russell [1872–1970] and the Italian mathematician Giuseppe Peano [1858–1932] than to that of Frege.) For an extensive discussion of logic symbols and operations, see LOGIC: *Logic systems: Formal logic*.

For some time, logicians were obsessed with a principle of parsimony, called Ockham’s razor, which justified them in reducing the number of these fundamental concepts, for example, by defining $p \supset q$ (read p implies q) as $\neg p \vee q$ or even as $\neg(p \wedge \neg q)$. While this definition, even if unnecessarily cumbersome, is legitimate classically, it is not permitted in intuitionistic logic (see below *Intuitionistic logic*). In the same spirit, many mathematicians adopted the Wiener-Kuratowski definition of the ordered pair $\langle a, b \rangle$ as $\{\{a\}, \{a, b\}\}$, where $\{a\}$ is the set whose sole element is a , which disguises its true significance.

Logic had been studied by the ancients, in particular by Aristotle and the Stoic philosophers. Philo of Megara (fl. c. 250 BC) had observed (or postulated) that $p \supset q$ is false

if and only if p is true and q is false. Yet the intimate connection between logic and mathematics had to await the insight of 19th-century thinkers, in particular Frege.

Frege was able to explain most mathematical notions with the help of his comprehension scheme, which asserts that, for every φ (formula or statement), there should exist a set X such that, for all x , $x \in X$ if and only if $\varphi(x)$ is true. Moreover, by the axiom of extensionality, this set X is uniquely determined by $\varphi(x)$. A flaw in Frege’s system was uncovered by Russell, who pointed out some obvious contradictions involving sets that contain themselves as elements—*e.g.*, by taking $\varphi(x)$ to be $\neg(x \in x)$. Russell illustrated this by what has come to be known as the barber paradox: A barber states that he shaves all who do not shave themselves. Who shaves the barber? Any answer contradicts the barber’s statement. To avoid these contradictions Russell introduced the concept of types, a hierarchy (not necessarily linear) of elements and sets such that definitions always proceed from more basic elements (sets) to more inclusive sets, hoping that self-referencing and circular definitions would then be excluded. With this type distinction, $x \in X$ only if X is of an appropriate higher type than x .

The type theory proposed by Russell, later developed in collaboration with the English mathematician Alfred North Whitehead (1861–1947) in their monumental *Principia Mathematica* (1910–13), turned out to be too cumbersome to appeal to mathematicians and logicians, who managed to avoid Russell’s paradox in other ways. Mathematicians made use of the Neumann-Gödel-Bernays set theory, which distinguishes between small sets and large classes, while logicians preferred an essentially equivalent first-order language, the Zermelo-Fraenkel axioms, which allow one to construct new sets only as subsets of given old sets. Mention should also be made of the system of the American philosopher Willard Van Orman Quine (b. 1908), which admits a universal set. (Cantor had not allowed such a “biggest” set, as the set of all its subsets would have to be still bigger.) Although type theory was greatly simplified by Alonzo Church and the American mathematician Leon Henkin (b. 1921), it came into its own only with the advent of category theory (see below *Category theory*).

Foundational logic. The prominence of logic in foundations led some people, referred to as logicians, to suggest that mathematics is a branch of logic. The concepts of membership and equality could reasonably be incorporated into logic, but what about the natural numbers? Kronecker had suggested that, while everything else was made by man, the natural numbers were given by God. The logicians, however, believed that the natural numbers were also man-made, inasmuch as definitions may be said to be of human origin.

Russell proposed that the number 2 be defined as the set of all two-element sets, that is, $X \in 2$ if and only if X has distinct elements x and y and every element of X is either x or y . The Hungarian-born American mathematician John von Neumann (1903–57) suggested an even simpler definition, namely that $X \in 2$ if and only if $X = 0$ or $X = 1$, where 0 is the empty set and 1 is the set consisting of 0 alone. Both definitions require an extralogical axiom to make them work—the axiom of infinity, which postulates the existence of an infinite set. Since the simplest infinite set is the set of natural numbers, one cannot really say that arithmetic has been reduced to logic. Most mathematicians follow Peano, who preferred to introduce the natural numbers directly by postulating the crucial properties of 0 and the successor operation S , among which one finds the principle of mathematical induction.

The logicist program might conceivably be saved by a 20th-century construction usually ascribed to Church, though he had been anticipated by the Austrian philosopher Ludwig Wittgenstein (1889–1951). According to Church, the number 2 is the process of iteration; that is, 2 is the function which to every function f assigns its iterate $2(f) = f \circ f$, where $(f \circ f)(x) = f(f(x))$. There are some theoretical difficulties with this construction, but these can be overcome if quantification over types is allowed; this is finding favour in theoretical computer science.

Einstein’s space-time manifold

Continuum hypothesis

Principia Mathematica

Impredicative constructions. A number of 19th-century mathematicians found fault with the program of reducing mathematics to arithmetic and set theory as suggested by the work of Cantor and Frege. In particular, the French mathematician Henri Poincaré (1854–1912) objected to impredicative constructions, which construct an entity of a certain type in terms of entities of the same or higher type—*i.e.*, self-referencing constructions and definitions. For example, when proving that every bounded nonempty set X of real numbers has a least upper bound a , one proceeds as follows. (For this purpose, it will be convenient to think of a real number, following Dedekind, as a set of rationals that contains all the rationals less than any element of the set.) One lets $x \in a$ if and only if $x \in y$ for some $y \in X$; but here y is of the same type as a .

It would seem that to do ordinary analysis one requires impredicative constructions. Russell and Whitehead tried unsuccessfully to base mathematics on a predicative type theory; but, though reluctant, they had to introduce an additional axiom, the axiom of reducibility, which rendered their enterprise impredicative after all. More recently, the Swedish logician Per Martin-Löf presented a new predicative type theory, but no one claims that this is adequate for all of classical analysis. However, the German-American mathematician Hermann Weyl (1885–1955) and the American mathematician Solomon Feferman have shown that impredicative arguments such as the above can often be circumvented and are not needed for most, if not all, of analysis. On the other hand, as was pointed out by the Italian computer scientist Giuseppe Longo (b. 1929), impredicative constructions are extremely useful in computer science—namely, for producing fixpoints (entities that remain unchanged under a given process).

Nonconstructive arguments. Another criticism of the Cantor-Frege program was raised by Kronecker, who objected to nonconstructive arguments, such as the following proof that there exist irrational numbers a and b such that a^b is rational. If $\sqrt{2}^{\sqrt{2}}$ is rational, then the proof is complete; otherwise take $a = \sqrt{2}^{\sqrt{2}}$ and $b = \sqrt{2}$, so that $a^b = 2$. The argument is nonconstructive, because it does not tell us which alternative holds, even though more powerful mathematics will, as was shown by the Russian mathematician Aleksandr Osipovich Gelfond (1906–68). In the present case, the result can be proved constructively by taking $a = \sqrt{2}$ and $b = 2 \log_2 3$. But there are other classical theorems for which no constructive proof exists.

Consider, for example, the statement

$$\exists_x (\exists_y \varphi(y) \supset \varphi(x)),$$

which symbolizes the statement that there exists a person who is famous if there are any famous people. This can be proved with the help of De Morgan's laws, named after the English mathematician and logician Augustus De Morgan (1806–71). It asserts the equivalence of $\exists_y \varphi(y)$ with $\neg \forall_y \neg \varphi(y)$, using classical logic, but there is no way one can construct such an x , for example, when $\varphi(x)$ asserts the existence of a well-ordering of the reals, as was proved by Feferman. An ordered set is said to be well-ordered if every nonempty subset has a least element. It had been shown by the German mathematician Ernst Zermelo (1871–1951) that every set can be well-ordered, provided one adopts another axiom, the axiom of choice, which says that, for every nonempty family of nonempty sets, there is a set obtainable by picking out exactly one element from each of these sets. This axiom is a fertile source of nonconstructive arguments.

Intuitionistic logic. The Dutch mathematician L.E.J. Brouwer (1881–1966) in the early 20th century had the fundamental insight that such nonconstructive arguments will be avoided if one abandons a principle of classical logic which lies behind De Morgan's laws. This is the principle of the excluded third (or excluded middle), which asserts that, for every proposition p , either p or not p ; and equivalently that, for every p , not not p implies p . This principle is basic to classical logic and had already been enunciated by Aristotle, though with some reservations, as he pointed out that the statement "there will be a sea battle tomorrow" is neither true nor false.

Brouwer did not claim that the principle of the excluded

third always fails, only that it may fail in the presence of infinite sets. Of two natural numbers x and y one can always decide whether $x = y$ or $x \neq y$, but of two real numbers this may not be possible, as one might have to know an infinite number of digits of their decimal expansions. Similar objections apply to De Morgan's laws, a consequence of the principle of the excluded third. For a finite set A , if it has been shown that the assertion $\forall_{x \in A} \neg \varphi(x)$ leads to a contradiction, $\exists_{x \in A} \varphi(x)$ can be verified by looking at each element of A in turn; *i.e.*, the statement that no members of a given set have a certain property can be disproved by examining in turn each element of the set. For an infinite set A , there is no way in which such an inspection can be carried out.

Brouwer's philosophy of mathematics is called intuitionism. Although Brouwer himself felt that mathematics was language-independent, his disciple Arend Heyting (1898–1980) set up a formal language for first-order intuitionistic arithmetic. Some of Brouwer's later followers even studied intuitionistic type theory (see below), which differs from classical type theory only by the absence of a single axiom (double negation):

$$\forall_{x \in \Omega} (\neg \neg x \supset x),$$

where Ω is the type of truth-values.

While it cannot be said that many practicing mathematicians have followed Brouwer in rejecting this principle on philosophical grounds, it came as a great surprise to people working in category theory that certain important categories called topoi (singular: topos; see below) have associated with them a language that is intuitionistic in general. In consequence of this fact, a theorem about sets proved constructively was immediately seen to be valid not only for sets but also for sheaves, which, however, lie beyond the scope of this article.

The moderate form of intuitionism considered here embraces Kronecker's constructivism but not the more extreme position of finitism. According to this view, which goes back to Aristotle, infinite sets do not exist, except potentially. In fact, it is precisely in the presence of infinite sets that intuitionists drop the classical principle of the excluded third.

An even more extreme position, called ultrafinitism, maintains that even very large numbers do not exist, say numbers greater than $10^{(10^{10})}$. Of course, the vast majority of mathematicians reject this view by referring to $10^{(10^{10})} + 1$, but the true believers have subtle ways of getting around this objection, which, however, lie beyond the scope of this discussion.

Other logics. While intuitionistic logic is obtained from classical logic by dropping the principle of the excluded third, other logics have also been proposed, though none has had a comparable impact on the foundations of mathematics. One may mention many-valued, or multivalued, logics, which admit a finite number of truth-values; fuzzy logic, with an imprecise membership relationship (though, paradoxically, a precise equality relation); and quantum logic, where conjunction may be only partially defined and implication may not be defined at all. Perhaps more important have been various so-called substructural logics in which the usual properties of the deduction symbol are weakened: relevance logic is studied by philosophers, linear logic by computer scientists, and a noncommutative version of the latter by linguists.

Formalism. Russell's discovery of a hidden contradiction in Frege's attempt to formalize set theory, with the help of his simple comprehension scheme, caused some mathematicians to wonder how one could make sure that no other contradictions existed. Hilbert's program, called formalism, was to concentrate on the formal language of mathematics and to study its syntax. In particular, the consistency of mathematics, which may be taken, for instance, to be the metamathematical assertion that the mathematical statement $0 = 1$ is not provable, ought to be a metathorem—that is, provable within the syntax of mathematics. This formalization project made sense only if the syntax of mathematics was consistent, for otherwise every syntactical statement would be provable, including that which asserts the consistency of mathematics.

Least upper bound

Well-ordered sets and the axiom of choice

Ultrafinitism

Unfortunately, a consequence of Gödel's incompleteness theorem (see below) is that the consistency of mathematics can be proved only in a language which is stronger than the language of mathematics itself. Yet, formalism is not dead—in fact, most pure mathematicians are tacit formalists—but the naive attempt to prove the consistency of mathematics in a weaker system had to be abandoned.

While no one, except an extremist intuitionist, will deny the importance of the language of mathematics, most mathematicians are also philosophical realists who believe that the words of this language denote entities in the real world. Following the Swiss mathematician Paul Bernays (1888–1977), this position is also called Platonism, since Plato believed that mathematical entities really exist.

Gödel. Implicit in Hilbert's program had been the hope that the syntactic notion of provability would capture the semantic notion of truth. Gödel came up with the surprising discovery that this was not the case for type theory and related languages adequate for arithmetic, as long as the following assumptions are insisted upon:

1. The set of theorems (provable statements) is effectively enumerable, by virtue of the notion of proof being decidable.
2. The set of true statements of mathematics is ω -complete in the following sense: given any formula $\varphi(x)$, containing a free variable x of type N , the universal statement $\forall_{x \in N} \varphi(x)$ will be true if $\varphi(n)$ is true for each numeral n —that is, for $n = 0, n = S0, n = SS0$, and so on.
3. The language is consistent.

Actually, Gödel also made a somewhat stronger assumption, which, as the American mathematician J. Barkley Rosser later showed, could be replaced by assuming consistency. Gödel's ingenious argument was based on the observation that syntactical statements about the language of mathematics can be translated into statements of arithmetic, hence into the language of mathematics. It was partly inspired by an argument that supposedly goes back to the ancient Greeks and which went something like this: Epimenides says that all Cretans are liars; Epimenides is a Cretan; hence Epimenides is a liar. Under the assumptions 1 and 2, Gödel constructed a mathematical statement g that is true but not provable. If it is assumed that all theorems are true, it follows that neither g nor $\neg g$ is a theorem.

No mathematician doubts assumption 1; by looking at a purported proof of a theorem, suitably formalized, it is possible for a mathematician, or even a computer, to tell whether it is a proof. By listing all proofs in, say, alphabetic order, an effective enumeration of all theorems is obtained. Classical mathematicians also accept assumption 2 and therefore reluctantly agree with Gödel that, contrary to Hilbert's expectation, there are true mathematical statements which are not provable.

However, moderate intuitionists could draw a different conclusion, because they are not committed to assumption 2. To them, the truth of the universal statement $\forall_{x \in N} \varphi(x)$ can be known only if the truth of $\varphi(n)$ is known, for each natural number n , in a uniform way. This would not be the case, for example, if the proof of $\varphi(n)$ increases in difficulty, hence in length, with n . Moderate intuitionists might therefore identify truth with provability and not be bothered by the fact that neither g nor $\neg g$ is true, as they would not believe in the principle of the excluded third in the first place.

Intuitionists have always believed that, for a statement to be true, its truth must be knowable. Moreover, moderate intuitionists might concede to formalists that to say that a statement is known to be true is to say that it has been proved. Still, some intuitionists do not accept the above argument. Claiming that mathematics is language-independent, intuitionists would state that in Gödel's metamathematical proof of his incompleteness theorem, citing ω -completeness to establish the truth of a universal statement yields a uniform proof of the latter after all.

Gödel considered himself to be a Platonist, inasmuch as he believed in a notion of absolute truth. He took it for granted, as do many mathematicians, that the set of true statements is ω -complete. Other logicians are more

skeptical and want to replace the notion of truth by that of truth in a model. In fact, Gödel himself, in his completeness theorem, had shown that for a mathematical statement to be provable it is necessary and sufficient that it be true in every model. His incompleteness theorem now showed that truth in every ω -complete model is not sufficient for provability. This point will be returned to later, as the notion of model for type theory is most easily formulated with the help of category theory, although this is not the way Gödel himself proceeded. See below *Gödel and category theory*.

Recursive definitions. Peano had observed that addition of natural numbers can be defined recursively thus:

$$x + 0 = x, x + Sy = S(x + y).$$

Other numerical functions $N^k \rightarrow N$ that can be defined with the help of such a recursion scheme (and with the help of 0, S, and substitution) are called primitive recursive. Gödel used this concept to make precise what he meant by "effectively enumerable." A set of natural numbers is said to be recursively enumerable if it consists of all $f(n)$ with $n \in N$, where $f: N \rightarrow N$ is a primitive recursive function.

This notion can easily be extended to subsets of N^k and, by a simple trick called arithmetization, to sets of strings of words in a language. Thus Gödel was able to assert that the set of theorems of mathematics is recursively enumerable, and, more recently, the American linguist Noam Chomsky (b. 1928) could say that the set of grammatical sentences of a natural language, such as English, is recursively enumerable.

It is not difficult to show that all primitive recursive functions can be calculated. For example, to calculate $x + y$ when $x = 3$ and $y = 2$, making use of Peano's recursive definition of $x + y$ and of the definitions $1 = S0, 2 = S1$, and so on, one proceeds as follows:

$$\begin{aligned} 3 + 2 &= S2 + S1 = S(S2 + 1) = S(S2 + S0) \\ &= SS(S2 + 0) = SSS2 = SS3 = S4 = 5. \end{aligned}$$

But primitive recursive functions are not the only numerical functions that can be calculated. More general are the recursive functions, where $f: N \rightarrow N$ is said to be recursive if its graph is recursively enumerable—that is, if there exist primitive recursive functions $u, v: N \rightarrow N$ such that, for all natural numbers x and y , $y = f(x)$ if and only if, for some $z \in N$, $x = u(z)$ and $y = v(z)$.

All recursive functions can be calculated with pencil and paper or, even more primitively, by moving pebbles (*calculi* in Latin) from one location to another, using some finite set of instructions, nowadays called a program. Conversely, only recursive functions can be so calculated, or computed by a theoretical machine introduced by the English mathematician Alan Turing (1912–54), or by a modern computer, for that matter. The Church-Turing thesis asserts that the informal notion of calculability is completely captured by the formal notion of recursive functions and hence, in theory, replicable by a machine.

Gödel's incompleteness theorem had proved that any useful formal mathematical system will contain undecidable propositions—propositions which can be neither proved nor disproved. Church and Turing, while seeking an algorithmic (mechanical) test for deciding theoremhood and thus potentially deleting nontheorems, proved independently, in 1936, that such an algorithmic method was impossible for the first-order predicate logic (see LOGIC, THE HISTORY AND KINDS OF: *20th-century logic*). The Church-Turing theorem of undecidability, combined with the related result of the Polish-born American mathematician Alfred Tarski (1902–83) on undecidability of truth, eliminated the possibility of a purely mechanical device replacing mathematicians.

Computers and proof. While many mathematicians use computers only as word processors and for the purpose of communication, computer-assisted computations can be useful for discovering potential theorems. For example, the prime number theorem was first suggested as the result of extensive hand calculations on the prime numbers up to 3,000,000 by the Swiss mathematician Leonhard Euler (1707–83), a process that would have been greatly facili-

Recursive enumerability

The Church-Turing thesis

The liar paradox

Absolute truth or truth within a model

Four-colour mapping theorem

tated by the availability of a modern computer. Computers may also be helpful in completing proofs when there are a large number of cases to be considered. The renowned computer-aided proof of the four-colour mapping theorem by the American mathematicians Kenneth Appel (b. 1932) and Wolfgang Haken (b. 1928) even goes beyond this, as the computer helped to determine which cases were to be considered in the next step of the proof. Yet, in principle, computers cannot be asked to discover proofs, except in very restricted areas of mathematics—such as elementary Euclidean geometry—where the set of theorems happens to be recursive, as was proved by Tarski.

As the result of earlier investigations by Turing, Church, the American mathematician Haskell Brooks Curry (1900–82), and others, computer science has itself become a branch of mathematics. Thus, in theoretical computer science, the objects of study are not just theorems but also their proofs, as well as calculations, programs, and algorithms. Theoretical computer science turns out to have a close connection to category theory. Although this lies beyond the scope of this article, an indication will be given below.

CATEGORY THEORY

Abstraction in mathematics. One recent tendency in the development of mathematics has been the gradual process of abstraction. The Norwegian mathematician Niels Henrik Abel (1802–29) proved that equations of the fifth degree cannot, in general, be solved by radicals. The French mathematician Évariste Galois (1811–32), motivated in part by Abel's work, introduced certain groups of permutations to determine the necessary conditions for a polynomial equation to be solvable. These concrete groups soon gave rise to abstract groups, which were described axiomatically. Then it was realized that to study groups it was necessary to look at the relation between different groups—in particular, at the homomorphisms which map one group into another while preserving the group operations. Thus people began to study what is now called the concrete category of groups, whose objects are groups and whose arrows are homomorphisms. It did not take long for concrete categories to be replaced by abstract categories, again described axiomatically. (For an introduction, see ALGEBRA: *Groups*.)

Categories

The important notion of a category was introduced by Samuel Eilenberg and Saunders Mac Lane at the end of World War II. These modern categories must be distinguished from Aristotle's categories, which are better called types in the present context. A category has not only objects but also arrows (referred to also as morphisms, transformations, or mappings) between them. For a more general introduction to categories, see ALGEBRA: *Categories*.

Many categories have as objects sets endowed with some structure and arrows, which preserve this structure. Thus, there exist the categories of sets (with empty structure) and mappings, of groups and group-homomorphisms, of rings and ring-homomorphisms, of vector spaces and linear transformations, of topological spaces and continuous mappings, and so on. There even exists, at a still more abstract level, the category of (small) categories and functors, as the morphisms between categories are called, which preserve relationships among the objects and arrows.

Functors

Not all categories can be viewed in this concrete way. For example, the formulas of a deductive system may be seen as objects of a category whose arrows $f: A \rightarrow B$ are deductions of B from A . In fact, this point of view is important in theoretical computer science, where formulas are thought of as types and deductions as operations.

More formally, a category consists of (1) a collection of objects A, B, C, \dots , (2) for each ordered pair of objects in the collection an associated collection of transformations including the identity $1_A: A \rightarrow A$, and (3) an associated law of composition for each ordered triple of objects in the category such that for $f: A \rightarrow B$ and $g: B \rightarrow C$ the composition gf (or $g \circ f$) is a transformation from A to C —i.e., $gf: A \rightarrow C$. Additionally, the associative law and the identities are required to hold (where the compositions are defined)—i.e., $h(gf) = (hg)f$ and $1_B f = f = f 1_A$.

In a sense, the objects of an abstract category have no windows, like the monads of Leibniz. To infer the interior of an object A one need only look at all the arrows from other objects to A . For example, in the category of sets, elements of a set A may be represented by arrows from a typical one-element set into A . Similarly, in the category of small categories, if $\mathbf{1}$ is the category with one object and no nonidentity arrows, the objects of a category A may be identified with the functors $\mathbf{1} \rightarrow A$. Moreover, if $\mathbf{2}$ is the category with two objects and one nonidentity arrow, the arrows of A may be identified with the functors $\mathbf{2} \rightarrow A$.

Isomorphic structures. An arrow $f: A \rightarrow B$ is called an isomorphism if there is an arrow $g: B \rightarrow A$ inverse to f —that is, such that $g \circ f = 1_A$ and $f \circ g = 1_B$. This is written $A \cong B$, and A and B are called isomorphic, meaning that they have essentially the same structure and that there is no need to distinguish between them. Inasmuch as mathematical entities are objects of categories, they are given only up to isomorphism. Their traditional set-theoretical constructions, aside from serving a useful purpose in showing consistency, are really irrelevant.

For example, in the usual construction of the ring of integers, an integer is defined as an equivalence class of pairs (m, n) of natural numbers, where (m, n) is equivalent to (m', n') if and only if $m + n' = m' + n$. The idea is that the equivalence class of (m, n) is to be viewed as $m - n$. What is important to a categorist, however, is that the ring \mathbf{Z} of integers is an initial object in the category of rings and homomorphisms—that is, that for every ring \mathbf{R} there is a unique homomorphism $\mathbf{Z} \rightarrow \mathbf{R}$. Seen in this way, \mathbf{Z} is given only up to isomorphism. In the same spirit, it should be said not that \mathbf{Z} is contained in the field \mathbf{Q} of rational numbers but only that the homomorphism $\mathbf{Z} \rightarrow \mathbf{Q}$ is one-to-one. Likewise, it makes no sense to speak of the set-theoretical intersection of π and $\sqrt{-1}$, if both are expressed as sets of sets of sets (ad infinitum).

Of special interest in foundations and elsewhere are adjoint functors (F, G) . These are pairs of functors between two categories \mathcal{A} and \mathcal{B} , which go in opposite directions such that a one-to-one correspondence exists between the set of arrows $F(A) \rightarrow B$ in \mathcal{B} and the set of arrows $A \rightarrow G(B)$ in \mathcal{A} —that is, such that the sets are isomorphic.

Adjoint functors

Topos theory. The original purpose of category theory had been to make precise certain technical notions of algebra and topology and to present crucial results of divergent mathematical fields in an elegant and uniform way, but it soon became clear that categories had an important role to play in the foundations of mathematics. This observation was largely the contribution of the American mathematician F.W. Lawvere (b. 1937), who elaborated on the seminal work of the German-born French mathematician Alexandre Grothendieck (b. 1928) in algebraic geometry. At one time he considered using the category of (small) categories (and functors) itself for the foundations of mathematics. Though he did not abandon this idea, later he proposed a generalization of the category of sets (and mappings) instead.

Among the properties of the category of sets, Lawvere singled out certain crucial ones, only two of which are mentioned here:

1. There is a one-to-one correspondence between subsets B of A and their characteristic functions $\chi: A \rightarrow \{\text{true}, \text{false}\}$, where, for each element a of A , $\chi(a) = \text{true}$ if and only if a is in B .
2. Given an element a of A and a function $h: A \rightarrow A$, there is a unique function $f: \mathbf{N} \rightarrow A$ such that $f(n) = h^n(a)$.

Suitably axiomatized, a category with these properties is called an (elementary) topos. However, in general, the two-element set $\{\text{true}, \text{false}\}$ must be replaced by an object Ω with more than two truth-values, though a distinguished arrow into Ω is still labeled as *true*.

Intuitionistic type theories. Topoi are closely related to intuitionistic type theories. Such a theory is equipped with certain types, terms, and theorems.

Among the types there should be a type Ω for truth-values, a type \mathbf{N} for natural numbers, and, for each type A , a type $\mathcal{P}(A)$ for all sets of entities of type A .

Among the terms there should be in particular:

1. The formulas $a = a'$ and $a \in \alpha$ of type Ω , if a and a' are of type A and α is of type $\mathcal{P}(A)$
2. The numerals 0 and Sn of type N , if the numeral n is of type N
3. The comprehension term $\{x \in A | \varphi(x)\}$ of type $\mathcal{P}(A)$, if $\varphi(x)$ is a formula of type Ω containing a free variable x of type A

The set of theorems should contain certain obvious axioms and be closed under certain obvious rules of inference, neither of which will be spelled out here.

At this point the reader may wonder what happened to the usual logical symbols. These can all be defined—for example, universal quantification

$$\forall_{x \in A} \varphi(x) \text{ as } \{x \in A | \varphi(x)\} = \{x \in A | x = x\}$$

and disjunction

$$p \vee q \text{ as } \forall_{t \in \Omega} ((p \supset t) \supset ((q \supset t) \supset t)).$$

For a formal definition of implication see LOGIC, THE HISTORY AND KINDS OF: *Logic systems: Formal logic*.

In general, the set of theorems will not be recursively enumerable. However, this will be the case for pure intuitionistic type theory \mathcal{L}_0 , in which types, terms, and theorems are all defined inductively. In \mathcal{L}_0 there are no types, terms, or theorems other than those that follow from the definition of type theory. \mathcal{L}_0 is adequate for the constructive part of the usual elementary mathematics—arithmetic and analysis—but not for metamathematics, if this is to include a proof of Gödel's completeness theorem, and not for category theory, if this is to include the Yoneda embedding of a small category into a set-valued functor category.

Internal language. It turns out that each topos \mathcal{T} has an internal language $L(\mathcal{T})$, an intuitionistic type theory whose types are objects and whose terms are arrows of \mathcal{T} . Conversely, every type theory \mathcal{L} generates a topos $T(\mathcal{L})$, by the device of turning (equivalence classes of) terms into objects, which may be thought of as denoting sets.

Nominalists may be pleased to note that every topos \mathcal{T} is equivalent (in the sense of category theory) to the topos generated by a language—namely, the internal language of \mathcal{T} . On the other hand, Platonists may observe that every type theory \mathcal{L} has a conservative extension to the internal language of a topos—namely, the topos generated by \mathcal{L} , assuming that this topos exists in the real (ideal) world. Here, the phrase “conservative extension” means that \mathcal{L} can be extended to $LT(\mathcal{L})$ without creating new theorems. The types of $LT(\mathcal{L})$ are names of sets in \mathcal{L} and the terms of $LT(\mathcal{L})$ may be identified with names of sets in \mathcal{L} for which it can be proved that they have exactly one element. This last observation provides a categorical version of Russell's theory of descriptions: if one can prove the unique existence of an x of type A in \mathcal{L} such that $\varphi(x)$, then this unique x has a name in $LT(\mathcal{L})$.

The interpretation of a type theory \mathcal{L} in a topos \mathcal{T} means an arrow $\mathcal{L} \rightarrow L(\mathcal{T})$ in the category of type theories or, equivalently, an arrow $T(\mathcal{L}) \rightarrow \mathcal{T}$ in the category of topoi. Indeed, T and L constitute a pair of adjoint functors.

Gödel and category theory. It is now possible to reexamine Gödel's theorems from a categorical point of view. In a sense, every interpretation of \mathcal{L} in a topos \mathcal{T} may be considered as a model of \mathcal{L} , but this notion of model is too general, for example, when compared with the models of classical type theories studied by Henkin. Therefore, it is preferable to restrict \mathcal{T} to being a special kind of topos called local. Given an arrow p into Ω in \mathcal{T} , then, p is true in \mathcal{T} if p coincides with the arrow true in \mathcal{T} , or, equivalently, if p is a theorem in the internal language of \mathcal{T} . \mathcal{T} is called a local topos provided that (1) $0 = 1$ is not true in \mathcal{T} , (2) $p \vee q$ is true in \mathcal{T} only if p is true in \mathcal{T} or q is true in \mathcal{T} , and (3) $\exists_{x \in A} \varphi(x)$ is true in \mathcal{T} only if $\varphi(a)$ is true in \mathcal{T} for some arrow $a: 1 \rightarrow A$ in \mathcal{T} . Here the statement $0 = 1$ in provision 1 can be replaced by any other contradiction—e.g., by $\forall_{t \in \Omega} t$, which says that every proposition is true.

A model of \mathcal{L} is an interpretation of \mathcal{L} in a local topos \mathcal{T} . Gödel's completeness theorem, generalized to intuitionistic type theory, may now be stated as follows: A closed formula of \mathcal{L} is a theorem if and only if it is true in every model of \mathcal{L} .

Gödel's incompleteness theorem, generalized likewise, says that, in the usual language of arithmetic, it is not enough to look only at ω -complete models: Assuming that \mathcal{L} is consistent and that the theorems of \mathcal{L} are recursively enumerable, with the help of a decidable notion of proof, there is a closed formula g in \mathcal{L} , which is true in every ω -complete model, yet g is not a theorem in \mathcal{L} .

The search for a distinguished model. A Platonist might still ask whether, among all the models of the language of mathematics, there is a distinguished model, which may be considered to be the world of mathematics. Take as the language \mathcal{L}_0 pure intuitionistic type theory (see above). It turns out, somewhat surprisingly, that the topos generated by \mathcal{L}_0 is a local topos; hence, the unique interpretation of \mathcal{L}_0 in the topos generated by it may serve as a distinguished model.

This so-called free topos has been constructed linguistically to satisfy any formalist, but it should also satisfy a moderate Platonist, one who is willing to abandon the principle of the excluded third, inasmuch as the free topos is the initial object in the category of all topoi. Hence, the free topos may be viewed, in the words of Leibniz, as the best of all possible worlds. More modestly speaking, the free topos is to an arbitrary topos like the ring of integers is to an arbitrary ring.

The language \mathcal{L}_0 should also satisfy any constructivist: if an existential statement $\exists_{x \in A} \varphi(x)$ can be proved in \mathcal{L}_0 , then $\varphi(a)$ can be proved for some term a of type A ; moreover, if $p \vee q$ can be proved, then either p can be proved or q can be proved.

The above argument would seem to make a strong case for the acceptance of pure intuitionistic type theory as the language of elementary mathematics—that is, of arithmetic and analysis—and hence for the acceptance of the free topos as the world of mathematics. Nonetheless, most practicing mathematicians prefer to stick to classical mathematics. In fact, classical arguments seem to be necessary for metamathematics—for example, in the usual proof of Gödel's completeness theorem—even for intuitionistic type theory.

In this connection, one celebrated consequence of Gödel's incompleteness theorem may be recalled, to wit: the consistency of \mathcal{L} cannot be proved (via arithmetization) within \mathcal{L} . This is not to say that it cannot be proved in a stronger metalanguage. Indeed, to exhibit a single model of \mathcal{L} would constitute such a proof.

It is more difficult to make a case for the classical world of mathematics, although this is what most mathematicians believe in. This ought to be a distinguished model of pure classical type theory \mathcal{L}_1 . Unfortunately, Gödel's argument shows that the interpretation of \mathcal{L}_1 in the topos generated by it is not a model in this sense.

Boolean local topoi. A topos is said to be Boolean if its internal language is classical. It is named after the English mathematician George Boole (1815–64), who was the first to give an algebraic presentation of the classical calculus of propositions. A Boolean topos is local under the following circumstances. The disjunction property (2) holds in a Boolean topos if and only if, for every closed formula p , either p is true or $\neg p$ is true. Moreover, with the help of De Morgan's laws, the existence property (3) may then be rephrased thus: if $\varphi(a)$ is true for all closed terms a of type A , then $\forall_{x \in A} \varphi(x)$ is true. As it turns out, a Boolean local topos may be described more simply, without referring to the internal language, as a topos with the following property: if $f, g: A \rightarrow B$ are arrows such that $fa = ga$ for all $a: 1 \rightarrow A$, then $f = g$. (Here 1 is the so-called terminal object, with the property that, for each object C , there is a unique arrow $C \rightarrow 1$.) For the Boolean topos to be ω -complete requires furthermore that all numerals—that is, closed terms of type N in its internal language—be standard—that is, have the form $0, S0, SS0$, and so on.

Of course, Gödel's completeness theorem shows that there are plenty of Boolean local topoi to model pure classical type theory in, but the usual proof of their existence requires nonconstructive arguments. It would be interesting to exhibit at least one such model constructively.

As a first step toward constructing a distinguished ω -complete Boolean model of \mathcal{L}_1 one might wish to de-

Pure intuitionistic type theory

Gödel's completeness theorems

fine the notion of truth in \mathcal{L}_1 , as induced by this model. Tarski had shown how truth can be defined for classical first-order arithmetic, a language that admits, aside from formulas, only terms of type N . Tarski achieved this essentially by incorporating ω -completeness into the definition of truth. It is not obvious whether his method can be extended to classical higher-order arithmetic—that is, to classical type theory. In fact, Tarski himself showed that the notion of truth is not definable (in a technical sense) in such a system. If his notion of definability corresponds to what is here meant by constructibility, then it is possible to conclude that, indeed, no Boolean model can be constructed.

The von Neumann universe

One may be tempted to consider as a candidate for the distinguished Boolean local topos the so-called von Neumann universe. This is defined as the union of a class of sets containing the empty set (the initial object in the category of sets) and closed under the power-set operation and under transfinite unions—thus, as a subcategory of the category of sets. But what is the category of sets if not the distinguished Boolean local topos being sought?

A better candidate may be Gödel's constructible universe, whose original purpose was to serve as a model of Zermelo-Fraenkel set theory in which the continuum hypothesis holds. It is formed like the von Neumann universe, except that the notion of subset, implicit in the power-set operation, is replaced by that of definable subset. Is it possible that this universe can be constructed syntactically, like the free topos, without reference to any previously given category of sets, or by a universal property?

Substitutional interpretation

In the internal language of a Boolean local topos, the logical connectives and quantifiers have their natural meanings. In particular, quantifiers admit a substitutional interpretation, a desirable property that has been discussed by philosophers (among them, Russell and the American logician Saul Kripke [b. 1941])—to wit: if an existential statement is true, then it can be witnessed by a term of appropriate type in the language; and a universal statement is true if it is witnessed by all terms of the appropriate type.

Note that, in the internal language of the free topos, and therefore in pure intuitionistic type theory, the substitutional interpretation is valid for existential quantifiers, by virtue of the free topos being local, but that it fails for universal quantifiers, in view of the absence of ω -completeness and the fact that in the free topos all numerals are standard. For a Boolean local topos, ω -completeness will also ensure that all numerals are standard, so that numerals mean exactly what they are intended to mean.

One distinguished model or many models. Some mathematicians do not believe that a distinguished world of mathematics should be sought at all, but rather that the multiplicity of such worlds should be looked at simultaneously. A major result in algebraic geometry, due to Alexandre Grothendieck, was the observation that every commutative ring may be viewed as a continuously variable local ring, as Lawvere would put it. In the same spirit, an amplified version of Gödel's completeness theorem would say that every topos may be viewed as a continuously variable local topos, provided sufficiently many

variables (Henkin constants) are adjoined to its internal language. Put in more technical language, this makes the possible worlds of mathematics stalks of a sheaf. However, the question still remains as to where this sheaf lives if not in a distinguished world of mathematics or—perhaps better to say—metamathematics.

These observations suggest that the foundations of mathematics have not achieved a definitive shape but are still evolving; they form the subject of a lively debate among a small group of interested mathematicians, logicians, and philosophers.

BIBLIOGRAPHY. W.S. ANGLIN and J. LAMBEK, *The Heritage of Thales* (1995), a textbook aimed primarily at undergraduate mathematics students, deals with the history, philosophy, and foundations of mathematics and includes an elementary introduction to category theory. Collections of important readings and original articles include PAUL BENACERRAF and HILARY PUTNAM (eds.), *Philosophy of Mathematics: Selected Readings*, 2nd ed. (1983), treating the foundations of mathematics, the existence of mathematical objects, the notion of mathematical truth, and the concept of set; JAAKO HINTIKKA (ed.), *The Philosophy of Mathematics* (1969), which includes articles by Henkin on completeness, by Feferman on predicativity, by Robinson on the calculus, and by Tarski on elementary geometry; and JEAN VAN HEIJENOORT (compiler), *From Frege to Gödel: A Source Book in Mathematical Logic, 1879–1931* (1967, reissued 1977). BERTRAND RUSSELL, *A History of Western Philosophy and Its Connection with Political and Social Circumstances from the Earliest Times to the Present Day*, 2nd ed. (1961, reprinted 1991), an extremely readable work, portrays the relevant views of the pre-Socratics, Plato, Aristotle, Leibniz, and Kant. MARIO BUNGE, *Treatise on Basic Philosophy*, vol. 7, *Epistemology & Methodology III. Philosophy of Science and Technology*, part 1, *Formal and Physical Sciences* (1985), contains a discussion by a philosopher of the different philosophical schools in the foundations of mathematics. WILLIAM KNEALE and MARTHA KNEALE, *The Development of Logic* (1962, reprinted 1984), offers a thorough scholarly account of the growth of logic from ancient times to the contributions by Frege, Russell, Brouwer, Hilbert, and Gödel. SAUNDERS MAC LANE, *Mathematics, Form and Function* (1986), records the author's personal views on the form and function of mathematics as a background to the philosophy of mathematics, touching on many branches of mathematics. MICHAEL HALLETT, *Cantorian Set Theory and Limitation of Size* (1984), provides a scholarly account of Cantor's set theory and its further development by Fraenkel, Zermelo, and von Neumann. WILLIAM S. HATCHER, *Foundations of Mathematics* (1968), surveys different systems, including those of Frege, of Russell, of von Neumann, Bernays, and Gödel, and of Quine as well as Lawvere's category of categories. Y.I. MANIN (I.U.I. MANIN), *A Course in Mathematical Logic*, trans. from Russian (1977), is addressed to mathematicians at a sophisticated level and presents the most significant discoveries up to 1977 concerning the continuum hypothesis, the nonexistence of algorithmic solutions, and other topics. GEORGE S. BOLOS and RICHARD C. JEFFREY, *Computability and Logic*, 3rd ed. (1989), for graduate and advanced undergraduate philosophy or mathematics students, deals with computability, Gödel's theorems, and the definability of truth, among other topics. J. LAMBEK and P.J. SCOTT, *Introduction to Higher Order Categorical Logic* (1986), is an advanced textbook addressed to graduate students in mathematics and computer science in which the relationship between topoi and type theories is explored in detail and some of the metatheorems cited in this article are proved. (J.L.)

The History of Mathematics

The following article is an account of how important parts of mathematics have developed historically.

As a consequence of the exponential growth of science most of this mathematics has developed since the 15th century AD, and it is a historical fact that from the 15th century to the late 20th century new developments in mathematics have been largely concentrated in Europe and North America. For these reasons the bulk of this article is devoted to European developments since 1500.

This does not mean, however, that developments elsewhere have been unimportant. Indeed, to understand the history of mathematics in Europe it is necessary to know its history at least in Mesopotamia and Egypt, in ancient Greece, and in Islāmic civilization from the 9th to the 15th centuries. The way in which these civilizations influenced one another, and the important direct contributions Greece and Islām made to later developments, are discussed in the first parts of this article.

India's contributions to the development of contemporary mathematics were made through the considerable influence of Indian achievements on Islāmic mathematics during its formative years. In order to provide a portrait of the mathematical achievements of one major Asian civilization, the article contains an overview of some of the principal periods and achievements of mathematics in China.

It is important to be aware of the character of the sources for the study of the history of mathematics. The history of Mesopotamian and Egyptian mathematics is based on the many extant original documents written by scribes. Although in the case of Egypt these documents are few, they are all of a type and leave little doubt that Egyptian mathematics was, on the whole, elementary and profoundly practical in its orientation. For Mesopotamian mathematics, on the other hand, there are a large number of clay tablets, which reveal mathematical achievements of a much higher order than those of the Egyptians. The tablets indicate that the Mesopotamians had a great deal of remarkable mathematical knowledge, although they offer no evidence that this knowledge was organized into a deductive system. Future research may reveal more about the early development of mathematics in Mesopotamia or about its influence on Greek mathematics, but it seems likely that this picture of Mesopotamian mathematics will stand.

From the period before Alexander the Great no Greek mathematical documents have been preserved except for fragmentary paraphrases, and even for the subsequent period it is well to remember that the oldest copies of Euclid's *Elements* are in Byzantine manuscripts dating from the 10th century AD. This stands in complete contrast to the situation described above for Egyptian and Babylonian documents. Although in general outline the present account of Greek mathematics is secure, in such important matters as the origin of the axiomatic method, the pre-Euclidean theory of ratios, and the discovery of the conic sections, historians have given competing accounts based on fragmentary texts, quotations of early writings culled from nonmathematical sources, and a considerable amount of conjecture.

Many important treatises from the early period of Islāmic mathematics have not survived or have survived only in Latin translations, so that there are still many unanswered questions about the relationship between early Islāmic mathematics and the mathematics of Greece and India. In addition, the amount of surviving material from later centuries is so large in comparison with that which has been studied that it is not yet possible to offer any sure judgment of what medieval Islāmic mathematics did not contain, and this means that it is not yet possible to evaluate with any assurance what was original in European mathematics from the 11th to the 15th century.

In modern times the invention of printing has largely solved the problem of obtaining secure texts and has allowed historians of mathematics to concentrate their editorial efforts on the correspondence or the unpublished works of mathematicians. However, the exponential growth of mathematics means that, for the period from the 19th century on, historians are able to treat only the major figures in any detail. In addition there is, as the period gets nearer the present, the problem of perspective. Mathematics, like any other human activity, has its fashions, and the nearer one is to a given period, the more likely these fashions are to look like the wave of the future. For this reason, the present article makes no attempt to assess the most recent developments in the subject.

(J.L.B.)

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, section 10/21, and the *Index*. The article is divided into the following sections:

-
- | | |
|--|---|
| Mathematics in ancient Mesopotamia 576 | European mathematics during the Middle Ages and Renaissance 587 |
| The numeral system and arithmetic operations | The transmission of Greek and Arabic learning |
| Geometric and algebraic problems | The universities |
| Mathematical astronomy | The Renaissance |
| Mathematics in ancient Egypt 577 | Mathematics in the 17th and 18th centuries 588 |
| The numeral system and arithmetic operations | The 17th century 588 |
| Geometry | Institutional background |
| Assessment of Egyptian mathematics | Numerical calculation |
| Greek mathematics 579 | Analytic geometry |
| The development of pure mathematics 579 | The calculus |
| The pre-Euclidean period | The 18th century 592 |
| The <i>Elements</i> | Institutional background |
| The three classical problems | Analysis and mechanics |
| Geometry in the 3rd century BC 581 | History of analysis |
| Archimedes | Other developments 594 |
| Apollonius | Theory of equations |
| Applied geometry | Foundations of geometry |
| Later trends in geometry and arithmetic 584 | Mathematics in the 19th and 20th centuries 595 |
| Greek trigonometry and mensuration | Projective geometry |
| Number theory | Making the calculus rigorous |
| Survival and influence of Greek mathematics | Fourier series |
| Mathematics in medieval Islām 585 | Elliptic functions |
| Origins | The theory of numbers |
| Mathematics in the 9th century | The theory of equations |
| Mathematics in the 10th century | Gauss |
| Omar Khayyam | Non-Euclidean geometry |
| Islāmic mathematics to the 15th century | |

Riemann
 Riemann's influence
 Differential equations
 Linear algebra
 The foundations of geometry
 The foundations of mathematics
 Cantor
 Mathematical physics
 Algebraic topology
 Developments in pure mathematics
 Mathematical physics and the theory of groups
 Mathematics in China and Japan 607

Chinese mathematics to the 13th century 607
 Outline of the history
 The *Nine Chapters*
 The commentary of Liu Hui
 The *Ten Classics of Mathematics*
 Some major developments from the 11th century
 to the 13th century
 The decline of the Sung-Yüan mathematics
 Japan in the 17th century 610
 The introduction of Chinese books
 The elaboration of Chinese methods
 Bibliography 610

Mathematics in ancient Mesopotamia

Until the 1920s it was commonly supposed that mathematics had its birth among the ancient Greeks. What was known of earlier traditions, such as the Egyptian as represented by the Rhind Papyrus (itself edited for the first time only in 1877), offered at best a meagre precedent. This impression gave way to a very different view as Orientalists succeeded in deciphering and interpreting the technical materials from ancient Mesopotamia.

Owing to the durability of the Mesopotamian scribes' clay tablets, the surviving evidence of this culture is substantial. Existing specimens of mathematics represent all the major eras—the Sumerian kingdoms of the 3rd millennium BC, the Akkadian and Babylonian regimes (2nd millennium), and the empires of the Assyrians (early 1st millennium), Persians (6th through 4th centuries BC), and Greeks (3rd century BC to 1st century AD). The level of competence was already high as early as the Old Babylonian dynasty, the time of the lawgiver king Hammurabi (c. 18th century BC), but after that there were few notable advances. The application of mathematics to astronomy, however, flourished during the Persian and Seleucid (Greek) periods.

The numeral system and arithmetic operations. Unlike the Egyptians, the mathematicians of the Old Babylonian period went far beyond the immediate challenges of their official accounting duties; for example, they introduced a versatile numeral system, which, like the modern system, exploited the notion of place value, and they developed computational methods that took advantage of this means of expressing numbers; they solved linear and quadratic problems by methods much like those now used in school algebra; their success with the study of what are now called Pythagorean number triples was a remarkable feat in number theory. The scribes who made such discoveries must have believed mathematics to be worthy of study in its own right, not just as a practical tool.

The older Sumerian system of numerals followed an additive decimal (base-10) principle similar to that of the Egyptians. But the Old Babylonian system converted this into a place-value system with the base of 60 (sexagesimal). The reasons for the choice of 60 are obscure, but one good mathematical reason might have been the existence of so many divisors (2, 3, 4, and 5, and some multiples) of the base, which would have greatly facilitated the operation of division. For numbers from 1 to 59, the symbols Υ for 1 and \triangleleft for 10 were combined in the simple additive manner (e.g., $\triangleleft\triangleleft\triangleleft\Upsilon\Upsilon$ represented 32). But, to express larger values, the Babylonians applied the concept of place value: for example, 60 was written as Υ , 70 as $\Upsilon\triangleleft$, 80 as $\Upsilon\triangleleft\triangleleft$, and so on. In fact, Υ could represent any power of 60. The context determined which power was intended. The Babylonians appear to have developed a placeholder symbol that functioned as a zero by the 3rd century BC, but its precise meaning and use is still uncertain. Furthermore, they had no mark to separate numbers into integral and fractional parts (as with the modern decimal point). Thus, the three-place numeral 3 7 30 could represent $3\frac{1}{8}$ (i.e., $3 + 7/60 + 30/60^2$), $187\frac{1}{2}$ (i.e., $3 \times 60 + 7 + 30/60$), 11,250 (i.e., $3 \times 60^2 + 7 \times 60 + 30$), or a multiple of these numbers by any power of 60.

The four arithmetic operations were performed in the same way as in the modern decimal system, except that carrying occurred whenever a sum reached 60 rather than 10. Multiplication was facilitated by means of tables; one

typical tablet lists the multiples of a number by 1, 2, 3, . . . , 19, 20, 30, 40, and 50. To multiply two numbers several places long, the scribe first broke the problem down into several multiplications, each by a one-place number, and then looked up the value of each product in the appropriate tables. He found the answer to the problem by adding up these intermediate results. These tables also assisted in division, for the values that head them were all reciprocals of regular numbers.

Regular numbers are those whose prime factors divide the base; the reciprocals of such numbers thus have only a finite number of places (by contrast, the reciprocals of nonregular numbers produce an infinitely repeating numeral). In base 10, for example, only numbers with factors of 2 and 5 (e.g., 8 or 50) are regular, and the reciprocals ($1/8 = 0.125$, $1/50 = 0.02$) have finite expressions; but the reciprocals of other numbers (such as 3 and 7) repeat infinitely ($0.\overline{3}$ and $0.14285\overline{7}$, respectively, where the bar indicates the digits that continually repeat). In base 60, only numbers with factors of 2, 3, and 5 are regular; for example, 6 and 54 are regular, so that their reciprocals ($1/6$ and $1/54$) are finite. The entries in the multiplication table for 1 6 40 are thus simultaneously multiples of its reciprocal $1/54$. To divide a number by any regular number, then, one can consult the table of multiples for its reciprocal.

An interesting tablet in the collection of Yale University (see Figure 1) shows a square with its diagonals; on one side is written "30," under one diagonal "42 25 35," and right along the same diagonal "1 24 51 10" (i.e., $1 + 24/60 + 51/60^2 + 10/60^3$). This third number is the correct value of $\sqrt{2}$ to four sexagesimal places (equivalent in the decimal system to 1.414213 . . . , which is too low by only 1 in the seventh place), while the second number is the product of the third number and the first and so gives the length of the diagonal when the side is 30. The scribe thus appears to have known an equivalent of the familiar long method of finding square roots. An additional element of sophistication is that, by choosing 30 (that is, $1/2$) for the side, the scribe obtained as the diagonal the reciprocal of

Yale Babylonian Collection



Figure 1: Babylonian mathematical tablet.

Multiplication

the value of $\sqrt{2}$ (since $\sqrt{2}/2 = 1/\sqrt{2}$), a result useful for purposes of division.

Geometric and algebraic problems. In a Babylonian tablet now in Berlin, the diagonal of a rectangle of sides 40 and 10 is solved as $40 + 10^2/(2 \times 40)$. Here, a very effective approximating rule is being used (that the square root of the sum of $a^2 + b^2$ can be estimated as $a + b^2/2a$), the same rule found frequently in later Greek geometric writings. Both these examples for roots illustrate the Babylonians' arithmetic approach in geometry. They also show that the Babylonians were aware of the relation between the hypotenuse and the two legs of a right triangle (now commonly known as the Pythagorean theorem) more than a thousand years before the Greeks used it.

A type of problem that occurs frequently in the Babylonian tablets seeks the base and height of a rectangle, where their product and sum have specified values. From the given information the scribe worked out the difference, since $(b - h)^2 = (b + h)^2 - 4bh$. In the same way, if the product and difference were given, the sum could be found. And, once both the sum and difference were known, each side could be determined, for $2b = (b + h) + (b - h)$ and $2h = (b + h) - (b - h)$. This procedure is equivalent to a solution of the general quadratic in one unknown. In some places, however, the Babylonian scribes solved quadratic problems in terms of a single unknown, just as would now be done by means of the quadratic formula.

Although these Babylonian quadratic procedures have often been described as the earliest appearance of algebra, there are important distinctions. The scribes lacked an algebraic symbolism; although they must certainly have understood that their solution procedures were general, they always presented them in terms of particular cases, rather than as the working through of general formulas and identities. They thus lacked the means for presenting general derivations and proofs of their solution procedures. Their use of sequential procedures rather than formulas, however, is less likely to detract from an evaluation of their effort now that algorithmic methods much like theirs have become commonplace through the development of computers.

As mentioned above, the Babylonian scribes knew that the base (b), height (h), and diagonal (d) of a rectangle satisfy the relation $b^2 + h^2 = d^2$. If one selects values at random for two of the terms, the third will usually be irrational, but it is possible to find cases in which all three terms are integers: for example, 3, 4, 5 and 5, 12, 13. (Such solutions are sometimes called Pythagorean triples.) A tablet in the Columbia University Collection presents a list of 15 such triples (decimal equivalents are shown in parentheses at the right; the gaps in the expressions for h , b , and d separate the place values in the sexagesimal numerals):

h	b	d			
2	1 59	2 49	(120	119	169)
57 36	56 7	1 20 45	(3,456	3,367	4,825)
1 20	1 16 41	1 50 49	(4,800	4,601	6,649)
3 45	3 31 49	5 9 1	(13,500	12,709	18,541)
1 12	1 5	1 37	(72	65	97)
...
1 30	56	1 46	(90	56	106)

(The entries in the column for h have to be computed from the values for b and d , for they do not appear on the tablet; but they must once have existed on a portion now missing.) The ordering of the lines becomes clear from another column, listing the values of d^2/h^2 (brackets indicate figures that are lost or illegible), which form a continually decreasing sequence: [1 59 0] 15, [1 56 56] 58 14 50 6 15, . . . , [1] 23 13 46 40. Accordingly, the angle formed between the diagonal and the base in this sequence increases continually from just over 45° to just under 60° . Other properties of the sequence suggest that the scribe knew the general procedure for finding all such number triples—that for any integers p and q , $2d/h = p/q + q/p$ and $2b/h = p/q - q/p$. (In the table, the implied values p and q turn out to be regular numbers falling in the standard set of reciprocals, as mentioned above in connection with the

multiplication tables.) Scholars are still debating nuances of the construction and the intended use of this table, but no one questions the high level of expertise implied by it.

Mathematical astronomy. The sexagesimal method developed by the Babylonians has a far greater computational potential than what was actually needed for the older problem texts. With the development of mathematical astronomy in the Seleucid period, however, it became indispensable. The astronomers sought to predict future occurrences of important phenomena, such as lunar eclipses and critical points in planetary cycles (conjunctions, oppositions, stationary points, and first and last visibility). They devised a technique for computing these positions (expressed in terms of degrees of latitude and longitude, measured relative to the path of the Sun's apparent annual motion) by successively adding appropriate terms in arithmetic progression. The results were then organized into a table listing positions as far ahead as the scribe chose. (Although the method is purely arithmetic, one can interpret it graphically: the tabulated values form a linear "zigzag" approximation to what is actually a sinusoidal variation.) While observations extending over centuries are required for finding the necessary parameters (e.g., periods, angular range between maximum and minimum values, and the like), only the computational apparatus at their disposal made the astronomers' forecasting effort possible.

Within a relatively short time (perhaps only a century or less), the elements of this system came into the hands of the Greeks. Although Hipparchus (2nd century BC) favoured the geometric approach of his Greek predecessors, he took over parameters from the Mesopotamians and adopted their sexagesimal style of computation. Through the Greeks it passed to Arabic scientists in the Middle Ages and thence to Europe, where it remained prominent in mathematical astronomy during the Renaissance and the early modern period. To this day it persists in the use of minutes and seconds to measure time and angles.

Aspects of the Old Babylonian mathematics may have come to the Greeks before this, early in the 5th century BC, the formative period of Greek geometry. There are a number of parallels that scholars have noted: for example, the Greek technique of "application of area" (see below) corresponded to the Babylonian quadratic methods (although in a geometric, not arithmetic, form). Further, the Babylonian rule for estimating square roots was widely used in Greek geometric computations, and there may also have been some shared nuances of technical terminology. Although details of the timing and manner of such a transmission are obscure because of the absence of explicit documentation, it seems that Western mathematics, while stemming largely from the Greeks, is considerably indebted to the older Mesopotamians.

Mathematics in ancient Egypt

The introduction of writing in Egypt in the predynastic period (c. 3000 BC) brought with it the formation of a special class of literate professionals, the scribes. By virtue of their writing skills, the scribes took on all the duties of a civil service: record keeping, tax accounting, the management of public works (building projects and the like), even the prosecution of war through overseeing military supplies and payrolls. Young men enrolled in scribal schools to learn the essentials of the trade, which included not only reading and writing but also the basics of mathematics.

One of the texts popular as a copy exercise in the schools of the New Kingdom (13th century BC) was a satirical letter in which one scribe, Hori, taunts his rival, Amenem-opet, for his incompetence as an adviser and manager. "You are the clever scribe at the head of the troops," Hori chides at one point; "a ramp is to be built, 730 cubits long, 55 cubits wide, with 120 compartments—it is 60 cubits high, 30 cubits in the middle . . . and the generals and the scribes turn to you and say, 'You are a clever scribe, your name is famous. Is there anything you don't know? Answer us, how many bricks are needed?' Let each compartment be 30 cubits by 7 cubits."

This problem, and three others like it in the same letter, cannot be solved without further data. But the point of

Calculation of planetary positions

Quadratic problems

Assessment of Egyptian mathematics. The papyri thus bear witness to a mathematical tradition closely tied to the practical accounting and surveying activities of the scribes. Occasionally, the scribes loosened up a bit: one problem (Rhind Papyrus, problem 79), for example, seeks the total from seven houses, seven cats per house, seven mice per cat, seven ears of wheat per mouse, and seven *hekat* of grain per ear (result: 19,607). The underlying scenario can only be guessed at, but it is probably of the playful “As I was going to St. Ives” type—certainly, the scribe’s interest in progressions (for which he appears to have a rule) goes beyond practical considerations. Other than this, however, Egyptian mathematics falls firmly within the range of practice.

Even allowing for the scantiness of the documentation that survives, the Egyptian achievement in mathematics must be viewed as modest. Its most striking features are competence and continuity. The scribes managed to work out the basic arithmetic and geometry necessary for their official duties as civil managers, and their methods persisted with little evident change for at least a millennium, perhaps two. Indeed, when Egypt came under Greek domination in the Hellenistic period (from the 3rd century BC onward), the older school methods continued. Quite remarkably, for example, the older unit-fraction methods are still prominent in Egyptian school papyri, written in the demotic (Egyptian) and Greek languages as late as the 7th century AD.

To the extent that Egyptian mathematics left a legacy at all, it was through its impact on the emerging Greek mathematical tradition between the 6th and 4th centuries BC. Because the documentation from this period is limited, the manner and significance of the influence can only be conjectured. But the report about Thales is only one of several such accounts of Greek intellectuals learning from Egyptians; Herodotus and Plato describe with approval Egyptian practices in the teaching and application of mathematics. This literary evidence has historical support, since the Greeks maintained continuous trade and military operations in Egypt from the 7th century BC onward. It is thus plausible that basic precedents for the Greeks’ earliest mathematical efforts—how they dealt with fractional parts or measured areas and volumes, or their use of ratios in connection with similar figures—came from the learning of the ancient Egyptian scribes.

Greek mathematics

THE DEVELOPMENT OF PURE MATHEMATICS

The pre-Euclidean period. The Greeks divided the field of mathematics into arithmetic (the study of “multitude,” or discrete quantity) and geometry (that of “magnitude,” or continuous quantity) and considered both to have originated in practical activities. Proclus, in his *Commentary on Euclid*, observes that geometry, literally, “measurement of land,” first arose in surveying practices among the ancient Egyptians, for the flooding of the Nile compelled them each year to redefine the boundaries of properties. Similarly, arithmetic started with the commerce and trade of Phoenician merchants. Although Proclus wrote quite late in the ancient period (in the 5th century AD), his account drew upon views proposed much earlier, by Herodotus (mid-5th century BC), for example, and by Eudemus, a disciple of Aristotle (late 4th century BC).

However plausible, this view is difficult to check, for there is only meagre evidence of practical mathematics from the early Greek period (roughly, the 8th through the 4th centuries BC). Inscriptions on stone, for example, reveal use of a numeral system the same in principle as the familiar Roman numerals. Herodotus seems to have known of the abacus as an aid for computation by both Greeks and Egyptians, and about a dozen stone specimens of Greek abaci survive from the 5th and 4th centuries BC. In the surveying of new cities in the Greek colonies of the 6th and 5th centuries, there was regular use of a standard length of 70 *plethra* (one *plethron* equals 100 feet) as the diagonal of a square of side 50 *plethra*; in fact, the actual diagonal of the square is $50\sqrt{2}$ *plethra*, so this was equivalent to using $\frac{7}{5}$ (or 1.4) as an estimate for

$\sqrt{2}$, which is now known to equal 1.414 . . . In the 6th century BC the engineer Eupalinus of Megara directed an aqueduct through a mountain on the island of Samos, and historians still debate how he did it. In a further indication of the practical aspects of early Greek mathematics, Plato describes in his *Laws* how the Egyptians drilled their children in practical problems in arithmetic and geometry; he clearly considered this a model for the Greeks to imitate.

Such hints about the nature of early Greek practical mathematics are confirmed in later sources, for example, in the arithmetic problems in papyrus texts from Ptolemaic Egypt (from the 3rd century BC onward) and the geometric manuals by Hero of Alexandria (1st century AD). In its basic manner this Greek tradition was much like the earlier traditions in Egypt and Mesopotamia. Indeed, it is likely that the Greeks borrowed from such older sources to some extent.

What was distinctive of the Greeks’ contribution to mathematics—and what in effect made them the creators of “mathematics,” as the term is usually understood—was its development as a theoretical discipline. This means two things: mathematical statements are general, and they are confirmed by proof. For example, the Mesopotamians had procedures for finding whole numbers a , b , and c for which $a^2 + b^2 = c^2$ (e.g., 3, 4, 5; 5, 12, 13; or 119, 120, 169). From the Greeks came a proof of a general rule for finding all such sets of numbers (now called Pythagorean triples): if one takes any whole numbers p and q , both being even or both odd, then $a = (p^2 - q^2)/2$, $b = pq$, and $c = (p^2 + q^2)/2$. As Euclid proves in Book X of the *Elements*, numbers of this form satisfy the relation for Pythagorean triples. Further, the Mesopotamians appear to have understood that sets of such numbers a , b , and c form the sides of right triangles, but the Greeks proved this result (Euclid, in fact, proves it twice, in *Elements*, Book I, proposition 47, and in a more general form in *Elements*, Book VI, proposition 31), and these proofs occur in the context of a systematic presentation of the properties of plane geometric figures.

The *Elements*, composed by Euclid of Alexandria, around 300 BC, was the pivotal contribution to theoretical geometry, but the transition from practical to theoretical mathematics had occurred much earlier, sometime in the 5th century BC. Initiated by men like Pythagoras of Samos (late 6th century) and Hippocrates of Chios (late 5th century), the theoretical form of geometry was advanced by others, most prominently the Pythagorean Archytas of Tarentum, Theaetetus of Athens, and Eudoxus of Cnidus (4th century). Because the actual writings of these men do not survive, knowledge about their work depends on remarks made by later writers. While even this limited evidence reveals how heavily Euclid depended on them, it does not set out clearly the motives behind their studies.

It is thus a matter of debate how and why this theoretical transition took place. A frequently cited factor is the discovery of the irrational. The early Pythagoreans held that “all things are number.” This might be taken to mean that any geometric measure can be associated with some number (that is, some whole number or fraction; in modern terminology, rational number), for in Greek usage the term for number, *arithmos*, refers exclusively to whole numbers or, in some contexts, to ordinary fractions. This assumption is common enough in practice, as when the length of a given line is said to be so many feet plus a fractional part. However, it breaks down for the lines that form the side and diagonal of the square. (For example, if it is supposed that the ratio between the side and diagonal may be expressed as the ratio of two whole numbers, it can be shown that both of these numbers must be even. This is impossible, since every fraction may be expressed as a ratio of two whole numbers having no common factors.) Geometrically, this means that there is no length that could serve as a unit of measure of both the side and diagonal; that is, the side and diagonal cannot each equal the same length multiplied by (different) whole numbers. Accordingly, the Greeks called such pairs of lengths “incommensurable.” (In modern terminology, unlike that of the Greeks, the term “number” is applied to such quantities as $\sqrt{2}$, but they are called irrational.)

Influence
on Greek
math-
ematics

Pythagorean
triples

The
theory of
irrationals

This result was already well known at the time of Plato and may well have been discovered within the school of Pythagoras in the 5th century BC, as some late authorities like Pappus of Alexandria (4th century AD) maintain. In any case, by 400 BC it was known that lines corresponding to $\sqrt{3}$, $\sqrt{5}$, and other square roots are incommensurable with a fixed unit length. The more general result, the geometric equivalent of the theorem that \sqrt{p} is irrational whenever p is not a rational square number, is associated with Plato's friend Theaetetus. Both Theaetetus and Eudoxus contributed to the further study of irrationals, and their followers collected the results into a substantial theory, as represented by the 115 propositions of Book X of the *Elements*.

The discovery of irrationals must have affected the very nature of early mathematical research, for it made clear that arithmetic was insufficient for the purposes of geometry, despite the assumptions made in practical work. Further, once such seemingly obvious assumptions as the commensurability of all lines turned out to be, in fact, false, then in principle all mathematical assumptions were rendered suspect. At the least, it became necessary to justify carefully all claims made about mathematics. Even more basically, it became necessary to establish what a reasoning has to be like to qualify as a proof. Apparently, Hippocrates of Chios, in the 5th century BC, and others soon after him had already begun the work of organizing geometric results into a systematic form in textbooks called "elements" (meaning "fundamental results" of geometry). These were to serve as sources for Euclid in his comprehensive textbook a century later.

The early mathematicians were not an isolated group but part of a larger, intensely competitive intellectual environment of pre-Socratic thinkers in Ionia and Italy, as well as Sophists at Athens. By insisting that only permanent things could have real existence, the philosopher Parmenides (5th century BC) called into question the most basic claims about knowledge itself. In contrast, Heraclitus (c. 500 BC) maintained that all permanence is an illusion, for the things that are perceived arise through a subtle balance of opposing tensions. What is meant by "knowledge" and "proof" thus came into debate.

Mathematical issues were often drawn into these debates. For some, like the Pythagoreans (and, later, Plato), the certainty of mathematics was held as a model for reasoning in other areas, like politics and ethics. But for others, mathematics seemed prone to contradiction. Zeno of Elea (5th century BC) posed paradoxes about quantity and motion. In one such paradox, it is assumed that a line can be bisected again and again without limit; if the division ultimately results in a set of points of zero length, then even infinitely many of them sum up only to zero, but, if it results in tiny line segments, then their sum will be infinite. In effect, the length of the given line must be both zero and infinite. In the 5th century BC a solution of such paradoxes was attempted by Democritus and the "atomists," philosophers who held that all material bodies are ultimately made up of invisibly small "atoms" (the Greek word *atomon* means "indivisible"). But in geometry such a view came into conflict with the existence of incommensurable lines, since the atoms would become the measuring units of all lines, even incommensurable ones. Protagoras and Democritus puzzled over whether the tangent to a circle meets it at a point or a line. The Sophists Antiphon and Bryson (both 5th century BC) considered how to compare the circle to polygons inscribed in it.

Influence
of the pre-
Socratics

The pre-Socratics thus revealed difficulties in specific assumptions about the infinitely many and the infinitely small, about the relation of geometry to physical reality, as well as in more general conceptions like "existence" and "proof." Philosophical questions such as these need not have affected the technical researches of mathematicians, but they did make them aware of difficulties that could bear on fundamental matters and so made them the more cautious in defining their subject matter.

Any such review of the possible effects of factors such as these is purely conjectural, since the sources are fragmentary and never make explicit how the mathematicians responded to the issues that were raised. But it is the par-

ticular concern over fundamental assumptions and proofs that distinguishes Greek mathematics from the earlier traditions. Plausible factors behind this concern can be identified in the special circumstances of the early Greek tradition—its technical discoveries and its cultural environment—even if it is not possible to describe in detail how these changes took place.

The Elements. The principal source for reconstructing pre-Euclidean mathematics is Euclid's *Elements*, for the major part of its contents can be traced back to research from the 4th century BC and in some cases even earlier. The first four books present constructions and proofs of plane geometric figures: Book I deals with the congruence of triangles, the properties of parallel lines, and the area relations of triangles and parallelograms; Book II establishes equalities relating to squares, rectangles, and triangles; Book III covers basic properties of circles; and Book IV sets out constructions of polygons in circles. Much of the content of Books I–III was already familiar to Hippocrates, and the material of Book IV can be associated with the Pythagoreans, so that this portion of the *Elements* has roots in 5th-century research. It is known, however, that questions about parallels were debated in Aristotle's school (around 350 BC), and so it may be assumed that efforts to prove results—such as the theorem stating that, for any given line and given point, there always exists a unique line through that point and parallel to the line—were tried and failed. Thus, the decision to found the theory of parallels on a postulate, as in Book I of the *Elements*, must have been a relatively recent development in Euclid's time. (The postulate would later become the subject of much study, and in modern times it led to the discovery of the so-called non-Euclidean geometries.)

Book V sets out a general theory of proportion, that is, a theory that does not require any restriction to commensurable magnitudes. This general theory derives from Eudoxus. On the basis of the theory, Book VI describes the properties of similar plane rectilinear figures and so generalizes the congruence theory of Book I. It appears that the technique of similar figures was already known in the 5th century BC, even though a fully valid justification could not have been given before Eudoxus worked out his theory of proportion.

Books VII–IX deal with what the Greeks called "arithmetic," the theory of whole numbers. It includes the properties of numerical proportions, greatest common divisors, least common multiples, and relative primes (Book VII); propositions on numerical progressions and square and cube numbers (Book VIII); and special results, like unique factorization into primes, the existence of an unlimited number of primes, and the formation of "perfect" numbers, that is, those numbers that equal the sum of their proper divisors (Book IX). In some form, Book VII stems from Theaetetus and Book VIII from Archytas.

The
theory of
numbers

Book X presents a theory of irrational lines and derives from the work of Theaetetus and Eudoxus. The remaining books treat the geometry of solids. Book XI sets out results on solid figures analogous to those for planes in Books I and VI; Book XII proves theorems on the ratios of circles, the ratios of spheres, and the volumes of pyramids and cones; Book XIII shows how to inscribe the five regular solids in a given sphere (compare the constructions of plane figures in Book IV). The measurement of curved figures in Book XII is inferred from that of rectilinear figures; for a particular curved figure a sequence of rectilinear figures is considered in which succeeding figures in the sequence become continually closer to the curved figure; the particular method used by Euclid derives from Eudoxus. The solid constructions in Book XIII derive from Theaetetus.

In sum the *Elements* gathered together the whole field of elementary geometry and arithmetic that had developed in the two centuries before Euclid. Doubtless, Euclid must be credited with particular aspects of this work, certainly with its editing as a comprehensive whole. But it is not possible to identify for certain even a single one of its results as having been his discovery. Other more advanced fields, though not touched on in the *Elements*, were already being vigorously studied in Euclid's time, in some

cases by Euclid himself. For these fields his textbook, true to its name, provides the appropriate “elementary” introduction.

One such field is the study of geometric constructions. Euclid, like geometers in the generation before him, divided mathematical propositions into two kinds: “theorems” and “problems.” A theorem makes the claim that all terms of a certain description have a specified property; a problem seeks the construction of a term that is to have a specified property. In the *Elements* all the problems are constructible on the basis of three stated postulates: that a line can be constructed by joining two given points, that a given line segment can be extended in a line indefinitely, and that a circle can be constructed with a given point as centre and a given line segment as radius. These postulates in effect restricted the constructions to the use of the so-called Euclidean tools—*i.e.*, a compass and a straightedge or unmarked ruler.

The Euclidean tools

The three classical problems. Although Euclid solves more than 100 construction problems in the *Elements*, many more were posed whose solutions required more than just compass and straightedge. Three such problems stimulated so much interest among later geometers that they have come to be known as the “classical problems”: doubling the cube, *i.e.*, constructing a cube whose volume is twice that of a given cube; trisecting the angle; and squaring the circle. Even in the pre-Euclidean period the effort to construct a square equal in area to a given circle had begun. Some related results came from Hippocrates, others were reported from Antiphon and Bryson, and Euclid’s theorem on the circle in *Elements*, Book XII, proposition 2, which states that circles are in the ratio of the squares of their diameters, was important for this search. But the first actual constructions (not, it must be noted, by means of the Euclidean tools, for this is impossible) came only in the 3rd century BC. The early history of angle trisection is obscure. Presumably it was attempted in the pre-Euclidean period, although solutions are known only from the 3rd century or later.

There are several successful efforts at doubling the cube, however, that date from the pre-Euclidean period. Hippocrates showed that the problem could be reduced to that of finding two mean proportionals: if for a given line a it is necessary to find x such that $x^3 = 2a^3$, lines x and y may be sought such that $a:x = x:y = y:2a$; for then $a^3/x^3 = (a/x)^3 = (a/x)(x/y)(y/2a) = a/2a = 1/2$. (Note that the same argument holds for any multiplier, not just the number 2.) Thus, the cube can be doubled if it is possible to find the two mean proportionals x and y between the two given lines a and $2a$. Constructions of the problem of the two means were proposed by Archytas, Eudoxus, and Menaechmus in the 4th century BC. Menaechmus, for example, constructed three curves corresponding to these same proportions: $x^2 = ay$, $y^2 = 2ax$, and $xy = 2a^2$; the intersection of any two of them then produces the line x that solves the problem. Menaechmus’ curves are conic sections: the first two are parabolas, the third a hyperbola. Thus, it is often claimed that Menaechmus originated the study of the conic sections. Indeed, Proclus and his older authority, Geminus (mid-1st century AD), appear to have held this view. The evidence does not indicate how Menaechmus actually conceived of the curves, however, so it is possible that the formal study of the conic sections as such did not begin until later, near the time of Euclid. Both Euclid and an older contemporary, Aristaeus, composed treatments (now lost) of the theory of conic sections.

In seeking the solutions of problems, geometers developed a special technique, which they called “analysis.” They assumed the problem to have been solved and then, by investigating the properties of this solution, worked back to find an equivalent problem that could be solved on the basis of the givens. To obtain the formally correct solution of the original problem, then, geometers reversed the procedure: first the data were used to solve the equivalent problem derived in the analysis, and from the solution obtained the original problem was then solved. In contrast to analysis, this reversed procedure is called “synthesis.”

Menaechmus’ cube duplication is an example of analysis: he assumed the mean proportionals x and y and then

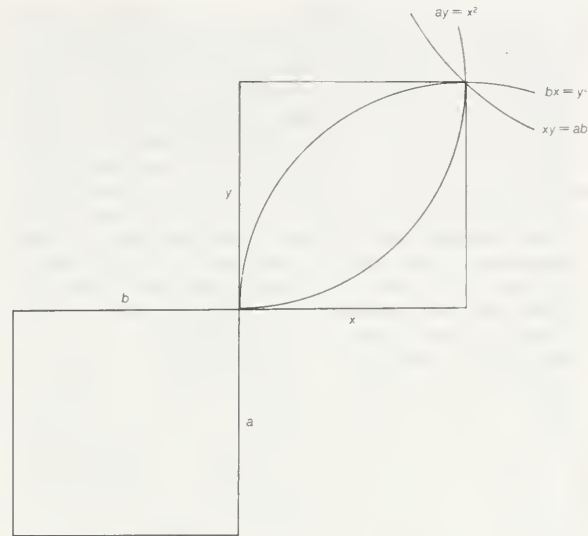


Figure 2: Cube duplication.

discovered them to be equivalent to the result of intersecting the three curves whose construction he could take as known. (The synthesis consists of introducing the curves, finding their intersection, and showing that this solves the problem.) It is clear that geometers of the 4th century BC were well acquainted with this method, but Euclid provides only syntheses, never analyses, of the problems solved in the *Elements*. Certainly in the cases of the more complicated constructions, however, there can be little doubt that some form of analysis preceded the syntheses presented in the *Elements*.

GEOMETRY IN THE 3RD CENTURY BC

The *Elements* was one of several major efforts by Euclid and others to consolidate the advances made over the 4th century BC. On the basis of these advances, Greek geometry entered its golden age in the 3rd century. This was a period rich with geometric discoveries, particularly in the solution of problems by analysis and other methods, and was dominated by the achievements of two figures: Archimedes of Syracuse (b. early 3rd century—d. 212 BC) and Apollonius of Perga (mid-3rd to early 2nd century BC).

Archimedes. Archimedes was most noted for his use of the Eudoxean method of exhaustion in the measurement of curved surfaces and volumes and for his applications of geometry to mechanics. To him is owed the first appearance and proof of the approximation $3\frac{1}{7}$ for the ratio of the circumference to the diameter of the circle (what is now designated π). Characteristically, Archimedes went beyond familiar notions, such as that of simple approximation, to more subtle insights, like the notion of bounds. For example, he showed that the perimeters of regular polygons circumscribed about the circle eventually become less than $3\frac{1}{7}$ the diameter as the number of their sides increases (Archimedes established the result for 96-sided polygons); similarly, the perimeters of the inscribed polygons eventually become greater than $3\frac{10}{71}$. Thus, these two values are upper and lower bounds, respectively, of π .

Archimedes’ result bears on the problem of circle quadrature in the light of another theorem he proved: that the area of a circle equals the area of a triangle whose height equals the radius of the circle and whose base equals its circumference. He established analogous results for the sphere showing that the volume of a sphere is equal to that of a cone whose height equals the radius of the sphere and whose base equals its surface area; the surface area of the sphere he found to be four times the area of its greatest circle. Equivalently, the volume of a sphere is shown to be two-thirds that of the cylinder which just contains it (that is, having height and diameter equal to the diameter of the sphere), while its surface is also equal to two-thirds that of the same cylinder (that is, if the circles that enclose the cylinder at top and bottom are included). The Greek historian Plutarch (early 2nd century AD) relates that Archimedes requested the figure for this theorem to

Volume of the sphere

Analysis and synthesis

be engraved on his tombstone, which is confirmed by the Roman writer Cicero (1st century BC), who actually located the tomb when he was quaestor of Sicily in 75 BC.

Apollonius. The work of Apollonius of Perga extended the field of geometric constructions far beyond the range in the *Elements*. For example, Euclid in Book III shows how to draw a circle so as to pass through three given points or to be tangent to three given lines; Apollonius (in a work called *Tangencies*, which no longer survives) found the circle tangent to three given circles, or tangent to any combination of three points, lines, and circles. (The three-circle tangency construction, one of the most extensively studied geometric problems, has attracted more than 100 different solutions in the modern period.)

Apollonius is best known for his *Conics*, a treatise in eight books (Books I–IV survive in Greek, V–VII in a medieval Arabic translation; Book VIII is lost). The conic sections are the curves formed when a plane intersects the surface of a cone (or double cone); it is assumed that the surface of the cone is generated by the rotation of a line through a fixed point around the circumference of a circle which is in a plane not containing that point. (The fixed point is the vertex of the cone, and the rotated line its generator.) There are three basic types: if the cutting plane is parallel to one of the positions of the generator, it produces the “parabola”; if it meets the cone only on one side of the vertex, it produces an “ellipse” (of which the circle is a special case), but, if it meets both parts of the cone, a “hyperbola.” Apollonius sets out in detail the properties of these curves. He shows, for example, that for given line segments a and b the parabola corresponds to the relation (in modern notation) $y^2 = ax$, the ellipse to $y^2 = ax - ax^2/b$, and the hyperbola to $y^2 = ax + ax^2/b$.

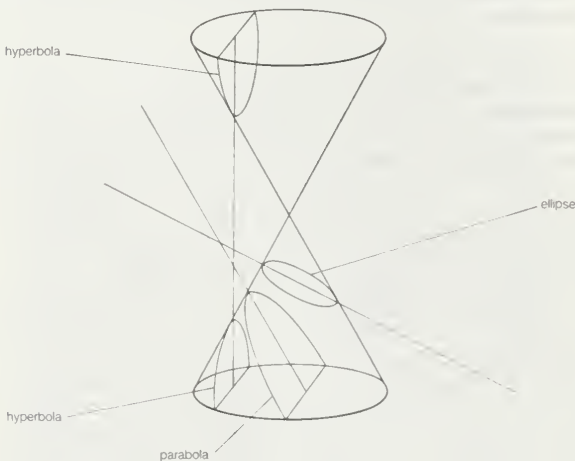


Figure 3: Conic sections.

Apollonius' treatise on conics in part consolidated more than a century of work before him and in part presented new findings of his own. As mentioned above, Euclid had already issued a textbook on the conics, while even earlier Menaechmus had played a role in their study. The names that Apollonius chose for the curves (the terms may be original with him) indicate yet an earlier connection. In the pre-Euclidean geometry *parabolē* referred to a specific operation, the “application” of a given area to a given line, in which the line x is sought such that $ax = b^2$ (where a and b are given lines); alternatively, x may be sought such that $x(a + x) = b^2$, or $x(a - x) = b^2$, and in these cases the application is said to be in “excess” (*hyperbolē*) or “defect” (*elleipsis*) by the amount of a square figure (namely, x^2). These constructions, which amount to a geometric solution of the general quadratic, appear in Books I, II, and VI of the *Elements* and can be associated in some form with the 5th-century Pythagoreans.

Apollonius presented a comprehensive survey of the properties of these curves. A sample of the topics he covered includes the following: the relations satisfied by the diameters and tangents of conics (Book I); how hyperbolas are related to their “asymptotes,” the lines they approach without ever meeting (Book II); how to draw tangents to

given conics (Book II); relations of chords intersecting in conics (Book III); the determination of the number of ways in which conics may intersect (Book IV); how to draw “normal” lines to conics, that is, lines meeting them at right angles (Book V); and the congruence and similarity of conics (Book VI).

By Apollonius' explicit statement, his results are of principal use as methods for the solution of geometric problems via conics. While he actually solved only a limited set of problems, the solutions of many others can be inferred from his theorems. For instance, the theorems of Book III permit the determination of conics that pass through given points or are tangent to given lines. In another work (now lost) Apollonius solved the problem of cube duplication by conics (a solution related in some way to that given by Menaechmus); further, a solution of the problem of angle trisection given by Pappus may have come from Apollonius or have been influenced by his work.

With the advance of the field of geometric problems by Euclid, Apollonius, and their followers, it became appropriate to introduce a classifying scheme: those problems solvable by means of conics were called solid, while those solvable by means of circles and lines only (as assumed in Euclid's *Elements*) were called planar. Thus, one can double the square by planar means (as in *Elements*, Book II, proposition 14), but one cannot double the cube in such a way, although a solid construction is possible (as given above). Similarly, the bisection of any angle is a planar construction (as shown in *Elements*, Book I, proposition 9), but the general trisection of the angle is of the solid type. It is not known when the classification was first introduced or when the planar methods were assigned canonical status relative to the others, but it seems plausible to date this near Apollonius' time. Indeed, much of his work—books like the *Tangencies*, the *Vergings* (or *Inclinations*), and the *Plane Loci*, now lost but amply described by Pappus—turns on the project of setting out the domain of planar constructions in relation to solutions by other means. On the basis of the principles of Greek geometry it cannot be demonstrated, however, that it is impossible to effect by planar means certain solid constructions (like the cube duplication and angle trisection). These results were only established by algebraists in the 19th century (notably by the French mathematician Pierre Laurent Wantzel in 1837).

A third class of problems, called linear, embraced those solvable by means of curves other than the circle and the conics (in Greek, the word for “line,” *grammē*, refers to all lines, whether curved or straight). For instance, one group of curves, the conchoids (from the Greek word for shell), are formed by marking off a certain length on a ruler and then pivoting it about a fixed point in such a way that one of the marked points stays on a given line; the other marked point traces out a conchoid. These curves can be used wherever a solution involves the positioning of a marked ruler relative to a given line (in Greek, such constructions are called *neuses*, or “vergings” of a line to a given point). For example, any acute angle (figured as

Linear problems

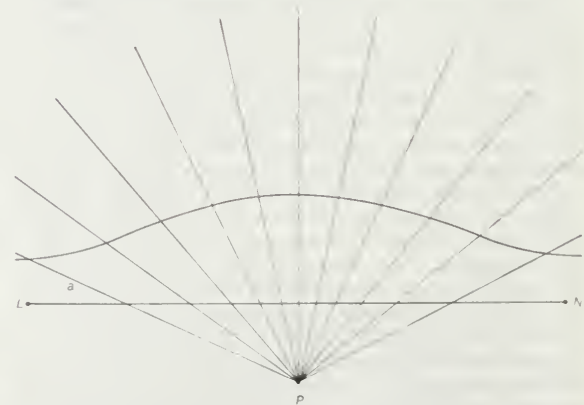


Figure 4: Conchoid curve. From fixed point P , several lines are drawn. A standard distance (a) is marked along each line from line LN , and the connection of the points creates a conchoid curve.

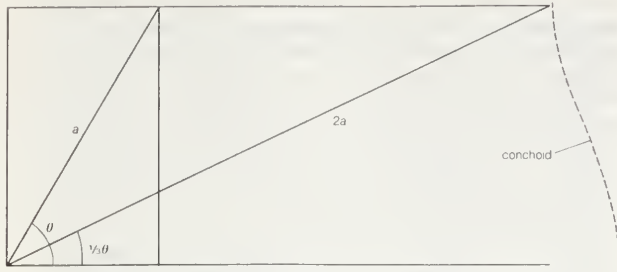


Figure 5: Angle trisection by means of a conchoid.

the angle between one side and the diagonal of a rectangle) can be trisected by taking a length equal to twice the diagonal and moving it about until it comes to be inserted between two other sides of the rectangle. If, instead, the appropriate conchoid relative to either of those sides is introduced, the required position of the line can be determined without the trial and error of a moving ruler. Because the same construction can be effected by means of a hyperbola, however, the problem is not linear but solid. Such uses of the conchoids were presented by Nicomedes (middle or late 3rd century BC), and their replacement by equivalent solid constructions appears to have come soon after, perhaps by Apollonius or his associates.

Some of the curves used for problem solving are not so reducible. For example, the Archimedean spiral couples uniform motion of a point on a half ray with uniform rotation of the ray around a fixed point at its end. Such curves have their principal interest as means for squaring the circle and trisecting the angle.

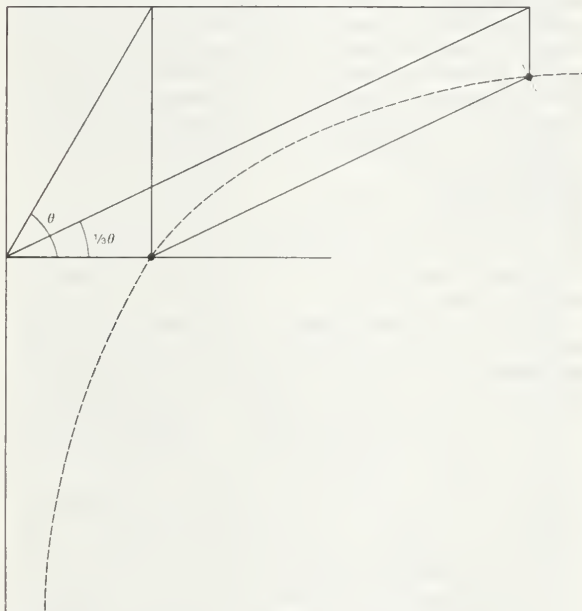


Figure 6: Angle trisection by means of a hyperbola.

Applied geometry. A major activity among geometers in the 3rd century BC was the development of geometric approaches in the study of the physical sciences, specifically, optics, mechanics, and astronomy. In each case the aim was to formulate the basic concepts and principles in terms of geometric and numerical quantities and then to derive the fundamental phenomena of the field by geometric constructions and proofs.

In optics, Euclid's textbook (called the *Optics*) set the precedent. Euclid postulated visual rays to be straight lines, and he defined the apparent size of an object in terms of the angle formed by the rays drawn from the top and the bottom of the object to the observer's eye. He then proved, for example, that nearer objects appear larger and appear to move faster and showed how to measure the height of distant objects from their shadows or reflected images, and so on. Other textbooks set out theorems on the phenomena of reflection and refraction (the field called catoptrics). The most extensive survey of optical phenomena is a treatise attributed to the astronomer

Ptolemy (2nd century AD), which survives only in the form of an incomplete Latin translation (12th century) based on a lost Arabic translation. It covers the fields of geometric optics and catoptrics, as well as experimental areas, such as binocular vision, and more general philosophical principles (the nature of light, vision, and colour). Of a somewhat different sort are the studies of burning mirrors by Diocles (late 2nd century BC), who proved that the surface that reflects the rays from the Sun to a single point is a paraboloid of revolution. Constructions of such devices remained of interest as late as the 6th century AD, when Anthemius of Tralles, best known for his work as architect of the Hagia Sophia at Constantinople, compiled a survey of remarkable mirror configurations.

Mechanics was dominated by the work of Archimedes, who was the first to prove the principle of balance: that two weights are in equilibrium when they are inversely proportional to their distances from the fulcrum. From this principle he developed a theory of the centres of gravity of plane and solid figures. He was also the first to state and prove the principle of buoyancy—that floating bodies displace their equal in weight—and to use it for proving the conditions of stability of segments of spheres and paraboloids, solids formed by rotating a parabolic segment about its axis. Archimedes proved the conditions under which this solid will return to its initial position if tipped, in particular, for the positions now called "stable I" and "stable II," where the vertex faces up and down, respectively.

In his work *Method Concerning Mechanical Theorems*, Archimedes also set out a special "mechanical method" that he used for the discovery of results on volumes and centres of gravity. He employed the bold notion of constituting solids from the plane figures formed as their sections (e.g., the circles that are the plane sections of spheres, cones, cylinders, and other solids of revolution), assigning to such figures a weight proportional to their area. For example, to measure the volume of a sphere, he imagined a balance beam, one of whose arms is a diameter of the sphere with the fulcrum at one endpoint of this diameter and the other arm an extension of the diameter to the other side of the fulcrum by a length equal to the diameter. Archimedes showed that the three circular cross sections made by a plane cutting the sphere and the associated cone and cylinder will be in balance (the circle in the cylinder with the circles in the sphere and cone) if the circle in the cylinder is kept in its original place, while the circles in the sphere and cone are placed with their centres of gravity at the opposite end of the balance. Doing this for all the sets of circles formed as cross sections of these solids by planes, he concluded that the solids themselves are in balance, the cylinder with the sphere and the cone together if the cylinder is left where it is, while the sphere and cone are placed with their centres of gravity at the opposite end of the balance. Since the centre of gravity of the cylinder is the midpoint of its axis, it follows that $(\text{sphere} + \text{cone}) : \text{cylinder} = 1 : 2$ (by the inverse proportion of weights and distances). Since the volume of the cone is one-third that of the cylinder, however, the volume of the sphere is found to be one-sixth that of the cylinder. In similar manner, Archimedes worked out the volumes and centres of gravity of spherical segments and segments of the solids of revolution of conic sections (paraboloids, ellipsoids, and hyperboloids). The critical notions—constituting solids out of their plane sections and assigning weights to geometric figures—were not formally valid within the standard conceptions of Greek geometry.

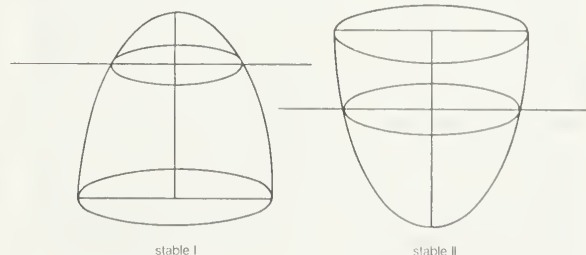


Figure 7: Paraboloids.

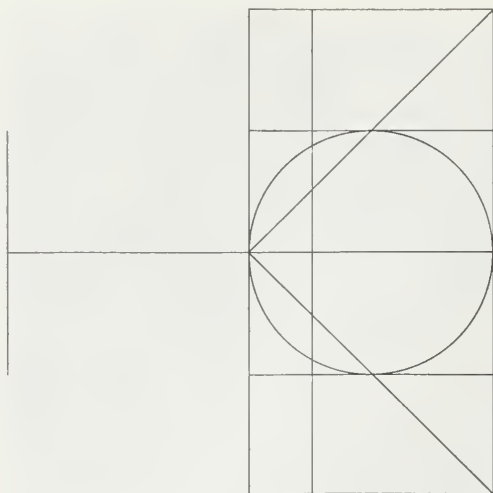


Figure 8: Measurement of the sphere.

and Archimedes admitted this. But he maintained that, although his arguments were not “demonstrations” (*i.e.*, proofs), they had value for the discovery of results about these figures.

Astronomy

The geometric study of astronomy has pre-Euclidean roots, with Eudoxus having developed a model for planetary motions around a stationary Earth. Accepting the principle—which, according to Eudemus, was first proposed by Plato—that only combinations of uniform circular motions are to be used, Eudoxus represented the path of a planet as the result of superimposing rotations of three or more concentric spheres whose axes are set at different angles. Although the fit with the phenomena was unsatisfactory, the curves thus generated (the *hippode*, or “horse-fetter”) continued to be of interest for their geometric properties, as is known through remarks by Proclus. Later geometers continued the search for geometric patterns satisfying the Platonic conditions. The simplest model, a scheme of circular orbits centred on the Sun, was introduced by Aristarchus of Samos (3rd century BC), but this was rejected by others, since a moving Earth was judged to be impossible on physical grounds. But Aristarchus’ scheme could have suggested use of an “eccentric” model, in which the planets rotate about the Sun and the Sun in turn rotates about the Earth. Apollonius introduced an alternative “epicyclic” model, in which the planet turns about a point that itself orbits in a circle (the “deferent”) centred at or near the Earth. As Apollonius knew, his epicyclic model is geometrically equivalent to an eccentric. These models were well adapted for explaining other phenomena of planetary motion. For instance, if the Earth is displaced from the centre of a circular orbit (as in the eccentric scheme), the orbiting body will appear to vary in speed (appearing faster when nearer the observer, slower when farther away), as is in fact observed for the Sun, Moon, and planets. By varying the relative sizes and rotation rates of the epicycle and deferent, in combination with the eccentric, a flexible device may be obtained for representing planetary motion.

LATER TRENDS IN GEOMETRY AND ARITHMETIC

Greek trigonometry and mensuration. After the 3rd century BC, mathematical research shifted increasingly away

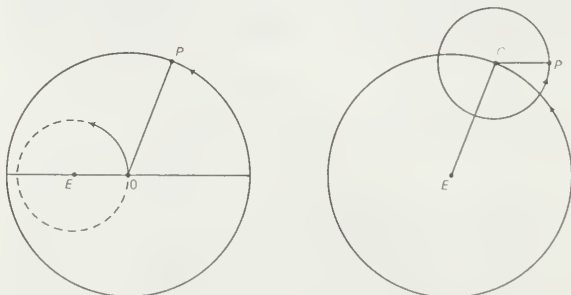


Figure 9: Eccentric and epicycle.

from the pure forms of constructive geometry toward areas related to the applied disciplines, in particular to astronomy. The necessary theorems on the geometry of the sphere (called spherics) were compiled into textbooks, such as the one by Theodosius (3rd or 2nd century BC), that consolidated the earlier work by Euclid and the work of Autolycus of Pitane (fl. c. 300 BC) on spherical astronomy. More significantly, in the 2nd century BC the Greeks first came into contact with the fully developed Mesopotamian astronomical systems and took from them many of their observations and parameters (for example, values for the average periods of astronomical phenomena). While retaining their own commitment to geometric models rather than adopting the arithmetic schemes of the Mesopotamians, the Greeks, nevertheless, followed the Mesopotamians’ lead in seeking a predictive astronomy based on a combination of mathematical theory and observational parameters. They thus made it their goal not merely to describe but to calculate the angular positions of the planets on the basis of the numerical and geometric content of the theory. This major restructuring of Greek astronomy, in both its theoretical and practical respects, was due primarily to Hipparchus (2nd century BC), whose work was consolidated and further advanced by Ptolemy.

To facilitate their astronomical researches, the Greeks developed techniques for the numerical measurement of angles, a precursor of trigonometry, and produced tables suitable for practical computation. Early efforts to measure the numerical ratios in triangles were made by Archimedes and Aristarchus. Their results were soon extended, and comprehensive treatises on the measurement of chords (in effect, a construction of a table of values equivalent to the trigonometric sine) were produced by Hipparchus and by Menelaus of Alexandria (early 2nd century AD). These works are now lost, but the essential theorems and tables are preserved in Ptolemy’s *Almagest* (Book 1, chapter 10). For computing with angles, the Greeks adopted the Mesopotamian sexagesimal method in arithmetic, whence it survives in the standard units for angles and time employed to this day.

Number theory. Although Euclid handed down a precedent for number theory in Books VII–IX of the *Elements*, later writers made no further effort to extend the field of theoretical arithmetic in his demonstrative manner. Beginning with Nicomachus of Gerasa (fl. c. AD 100), several writers produced collections expounding a much simpler form of number theory. A favourite result is the representation of arithmetic progressions in the form of “polygonal numbers.” For instance, if the numbers 1, 2, 3, 4, . . . are added successively, the “triangular” numbers 1, 3, 6, 10, . . . are obtained; similarly, the odd numbers 1, 3, 5, 7, . . . sum to the “square” numbers 1, 4, 9, 16, . . . , while the sequence 1, 4, 7, 10, . . . with a constant difference of 3, sums to the “pentagonal” numbers 1, 5, 12, 22, In general, these results can be expressed in the form of geometric shapes formed by lining up dots in the appropriate two-dimensional configurations. In the ancient arithmetics, such results are invariably presented as particular cases, without any general notational method or general proof. The writers in this tradition are called neo-Pythagoreans, since they viewed themselves as continuing the Pythagorean school of the 5th century BC, and in the spirit of ancient Pythagoreanism they tied their numerical interests to a philosophical theory that was an amalgam of Platonic metaphysical and theological doctrines. With its exponent Iamblichus of Chalcis (4th century AD), neo-Pythagoreans became a prominent part of the revival of pagan religion in opposition to Christianity in late antiquity.

An interesting concept of this school of thought, which Iamblichus attributes to Pythagoras himself, is that of “amicable numbers”: two numbers are amicable if each is equal to the sum of the proper divisors of the other (for example, 220 and 284). Attributing virtues such as friendship and justice to numbers was characteristic of the Pythagoreans at all times.

Of much greater mathematical significance is the arithmetic work of Diophantus of Alexandria (active at an unknown time between the 2nd century BC and the 3rd

The measurement of angles

Diophantus of Alexandria

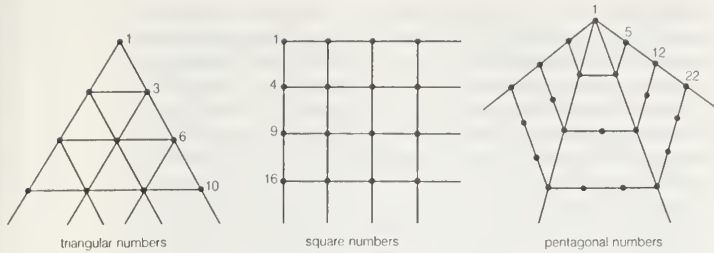


Figure 10: Polygonal arrays.

century AD). His writing, the *Arithmetica*, originally in 13 books (six survive in Greek, another four in the medieval Arabic translation), sets out hundreds of arithmetic problems with their solutions. For example, Book II, problem 8, seeks to express a given square number as the sum of two square numbers (here and throughout, the “numbers” are rational). Like those of the neo-Pythagoreans, his treatments are always of particular cases rather than general solutions; thus, in this problem the given number is taken to be 16 and the solutions worked out are $2^{56}/25$ and $1^{44}/25$. In this example, as is often the case, the solutions are not unique; indeed, in the very next problem Diophantus shows how a number given as the sum of two squares (e.g., $13 = 4 + 9$) can be expressed differently as the sum of two other squares (for example, $13 = 3^{24}/25 + 1/25$).

To find his solutions, Diophantus adopted an arithmetic form of the method of analysis. He first reformulated the problem in terms of one of the unknowns, and he then manipulated it as if it were known until an explicit value for the unknown emerged. He even adopted an abbreviated notational scheme to facilitate such operations, where, for example, the unknown is symbolized by a figure somewhat resembling the Roman letter *S*. (This is a standard abbreviation for the word number in ancient Greek manuscripts.) Thus, in the first problem discussed above, if *S* is one of the unknown solutions, then $16 - S^2$ is a square, supposing the other unknown to be $2S - 4$ (where the 2 is arbitrary but the 4 chosen because it is the square root of the given number 16), Diophantus found from summing the two unknowns ($[2S - 4]^2$ and S^2) that $4S^2 - 16S + 16 + S^2 = 16$, or $5S^2 = 16S$, that is, $S = 16/5$. So one solution is $S^2 = 2^{56}/25$, while the other solution is $16 - S^2$, or $1^{44}/25$.

Survival and influence of Greek mathematics. Notable in the closing phase of Greek mathematics were Pappus (early 4th century AD), and Theon (late 4th century) and his daughter Hypatia (d. 415). All were active in Alexandria as professors of mathematics and astronomy, and they produced extensive commentaries on the major authorities—Pappus and Theon on Ptolemy, Hypatia on Diophantus and Apollonius. Later, Eutocius of Ascalon (early 6th century) produced commentaries on Archimedes and Apollonius. While much of their output has since been lost, much survives. They proved themselves reasonably competent in technical matters but little inclined toward significant insights (their aim was usually to fill in minor steps assumed in the proofs, to append alternative proofs, and the like), and the level of originality was very low. But these scholars frequently preserved fragments of older works that are now lost, and their teaching and editorial efforts assured the survival of the works of Euclid, Archimedes, Apollonius, Diophantus, Ptolemy, and others that now do exist, either in Greek manuscripts or in medieval translations (Arabic, Hebrew, and Latin) derived from them.

The legacy of Greek mathematics, particularly in the fields of geometry and geometric science, was enormous. From an early period the Greeks formulated the objectives of mathematics not in terms of practical procedures but as a theoretical discipline committed to the development of general propositions and formal demonstrations. The range and diversity of their findings, especially those of the masters of the 3rd century BC, supplied geometers with subject matter for centuries thereafter, even though the tradition that was transmitted into the Middle Ages and Renaissance was incomplete and defective.

The rapid rise of mathematics in the 17th century was

based in part on the conscious imitation of the ancient classics and on competition with them. In the geometric mechanics of Galileo and the infinitesimal researches of Kepler and Cavalieri, it is possible to perceive a direct inspiration from Archimedes. The study of the advanced geometry of Apollonius and Pappus stimulated new approaches in geometry—for example, the analytic methods of Descartes and the projective theory of Girard Desargues. Purists like Huygens and Newton insisted on the Greek geometric style as a model of rigour, just as others sought to escape its forbidding demands of completely worked-out proofs. The full impact of Diophantus’ work is evident particularly with Pierre de Fermat in his researches in algebra and number theory. Although mathematics has today gone far beyond the ancient achievements, the leading figures of antiquity, like Archimedes, Apollonius, and Ptolemy, can still be rewarding reading for the ingenuity of their insights. (W.R.K.)

Mathematics in medieval Islām

Origins. In Hellenistic times and in late antiquity, scientific learning in the eastern part of the Roman world was spread over a variety of centres, and Justinian’s closing of the pagan academies in Athens in 529 gave further impetus to this diffusion. An additional factor was the translation and study of Greek scientific and philosophical texts sponsored both by monastic centres of the various Christian churches in the Levant, Egypt, and Mesopotamia and by enlightened rulers of the Sasanian dynasty in places like the medical school at Gondeshapur.

Also important were developments in India in the first few centuries AD. Although the decimal system for whole numbers was apparently not known to the Indian astronomer Āryabhaṭa I (b. 476), it was used by his pupil Bhāskara I in 620, and by 670 the system had reached northern Mesopotamia, where the Nestorian bishop Severus Sebokht praised its Hindu inventors as discoverers of things more ingenious than those of the Greeks. Earlier, in the late 4th or early 5th century, the anonymous Hindu author of an astronomical handbook, the *Sūrya Siddhānta*, had tabulated the sine function (unknown in Greece) for every $3\frac{3}{4}^\circ$ of arc from $3\frac{3}{4}^\circ$ to 90° .

Within this intellectual context the rapid expansion of Islām took place between the time of Muḥammad’s return to Mecca from his exile in Medina in 630 and the Muslim conquest of lands extending from Spain to the borders of China by 715. Not long afterward, Muslims began the acquisition of foreign learning, and, by the time of the caliph al-Manṣūr (d. 775), such Indian and Persian astronomical material as the *Brāhma-sphuṭa-siddhānta* and the *Shāh’s Tables* had been translated into Arabic. The subsequent acquisition of Greek material was greatly advanced when the caliph al-Ma’mūn constructed a translation and research centre, the House of Wisdom, in Baghdad during his reign (813–833). Most of the translations were done from Greek and Syriac by Christian scholars, but the impetus and support for this activity came from Muslim patrons. These included not only the caliph but also wealthy individuals such as the three brothers known as the Banū Mūsā, whose treatises on geometry and mechanics formed an important part of the works studied in the Islāmic world.

Of Euclid’s works, the *Elements*, the *Data*, the *Optics*, the *Phaenomena*, and *On Divisions* were translated. Of Archimedes’ works only two—*Sphere and Cylinder* and *Measurement of the Circle*—are known to have been translated, but these were sufficient to stimulate independent researches from the 9th to the 15th century. On the other hand, virtually all of Apollonius’ works were translated, and of Diophantus and Menelaus one book each, the *Arithmetica* and the *Sphaerica*, respectively, were translated into Arabic. Finally, the translation of Ptolemy’s *Almagest* furnished important astronomical material.

Of the minor writings, Dioctes’ treatise on mirrors, Theodosius’ *Spherics*, Pappus’ work on mechanics, Ptolemy’s *Planisphaerium*, and Hypsicles’ treatises on regular polyhedra (the so-called Books XIV and XV of Euclid’s *Elements*) were among those translated.

Indian astronomy and mathematics

Translation of Greek texts

Mathematics in the 9th century. Thābit ibn Qurrah (836–901), a Sabian from Ḥarrān in northern Mesopotamia, was an important translator and reviser of these Greek works. In addition to translating works of the major Greek mathematicians (for the Banū Mūsā, among others), he was a court physician. He also translated Nicomachus of Gerasa's *Arithmetic* and discovered a beautiful rule for finding amicable numbers, a pair of numbers such that each number is the sum of the set of proper divisors of the other number. The investigation of such numbers formed a continuing tradition in Islām. Kamāl ad-Dīn al-Fārisī (d. c. 1320) gave the pair 17,926 and 18,416 as an example of Thābit's rule, and in the 17th century Muḥammad Bāqir Yazdī gave the pair 9,363,584 and 9,437,056.

One scientist typical of the 9th century was Muḥammad ibn Mūsā al-Khwārizmī. Working in the House of Wisdom, he introduced Indian material in his astronomical works and also wrote an early book explaining Hindu arithmetic, the *Book of Addition and Subtraction According to the Hindu Calculation*. In another work, the *Book of Restoring and Balancing*, he provided a systematic introduction to algebra, including a theory of quadratic equations. Both works had important consequences for Islāmic mathematics. *Hindu Calculation* began a tradition of arithmetic books that, by the middle of the next century, led to the invention of decimal fractions (complete with a decimal point), and his *Restoring and Balancing* became the point of departure and model for later writers such as the Egyptian Abū Kāmil. Both books were translated into Latin, and *Restoring and Balancing* was the origin of the word algebra, from the Arabic word for "restoring" in its title (*al-jabr*). The *Hindu Calculation*, from a Latin form of the author's name, *algorismi*, yielded the word algorithm.

Al-Khwārizmī's algebra also served as a model for later writers in its application of arithmetic and algebra to the distribution of inheritances according to the complex requirements of Muslim religious law. This tradition of service to the Islāmic faith was an enduring feature of mathematical work in Islām and one which, in the eyes of many, justified the study of secular learning. In the same category are al-Khwārizmī's method of calculating the time of visibility of the new moon (which signals the beginning of the Muslim month) and the expositions by astronomers of methods for finding the direction to Mecca for the five daily prayers.

Mathematics in the 10th century. Islāmic scientists in the 10th century were involved in three major mathematical projects: the completion of arithmetic algorithms, the development of algebra, and the extension of geometry.

The first of these projects led to the appearance of three complete numeration systems, one of which was the finger arithmetic used by the scribes and treasury officials. This ancient arithmetic system, which became known throughout the East and Europe, employed mental arithmetic and a system of storing intermediate results on the fingers as an aid to memory. (Its use of unit fractions recalls the Egyptian system.) During the 10th and 11th centuries capable mathematicians, such as Abū al-Wafā' (940–997/8), wrote on this system, but it was eventually replaced by the decimal system.

A second common system was the base-60 numeration inherited from the Babylonians via the Greeks and known as the arithmetic of the astronomers. Although astronomers used this system for their tables, they usually converted numbers to the decimal system for complicated calculations and then converted the answer back to sexagesimals.

The third system was Indian arithmetic, whose basic numeral forms, complete with the zero, eastern Islām took over from the Hindus. (Different forms of the numerals, whose origins are not entirely clear, were used in the western part of Islām.) The basic algorithms also came from India, but these were adapted by al-Uqlidīsī (c. 950) to pen and paper instead of the traditional dust board, a move that helped to popularize this system. Also, the arithmetic algorithms were completed in two ways: by the extension of root-extraction procedures, known to Hindus

and Greeks only for square and cube roots, to roots of higher degree, and by the extension of the Hindu decimal system for whole numbers to include decimal fractions. These fractions appear simply as computational devices in the work of both al-Uqlidīsī and al-Baghdādī (c. 1000), but in subsequent centuries they received systematic treatment as a general method. As for extraction of roots, Abū al-Wafā' wrote a treatise (now lost) on the topic, and Omar Khayyam ('Umar al-Khayyāmī [1048–1131]) solved the general problem of extracting roots of any desired degree. Omar's treatise, too, is lost, but the method is known from other writers, and it appears that a major step in its development was al-Karajī's 10th-century derivation by means of mathematical induction of the binomial theorem for whole-number exponents—i.e., his discovery that

$$(a + b)^n = a^n + na^{n-1}b + \frac{n(n-1)}{2} a^{n-2}b^2 + \frac{n(n-1)(n-2)}{2 \cdot 3} a^{n-3}b^3 + \dots + nab^{n-1} + b^n.$$

During the 10th century Islāmic algebraists progressed from al-Khwārizmī's quadratic polynomials to the mastery of the algebra of expressions involving arbitrary positive or negative integral powers of the unknown. Several algebraists explicitly stressed the analogy between the rules for working with powers of the unknown in algebra and those for working with powers of 10 in arithmetic, and there was interaction between the development of arithmetic and algebra from the 10th to the 12th century. A 12th-century student of al-Karajī's works, as-Samaw'al, was able to approximate the quotient $(20x^2 + 30x)/(6x^2 + 12)$ as

$$3\frac{1}{3} + 5\left(\frac{1}{x}\right) - 6\frac{2}{3}\left(\frac{1}{x^2}\right) - 10\left(\frac{1}{x^3}\right) + \dots - 40\left(\frac{1}{x^7}\right)$$

and also gave a rule for finding the coefficients of the successive powers of $1/x$. Although none of this employed symbolic algebra, algebraic symbolism was in use by the 14th century in the western part of the Islāmic world. The context for this well-developed symbolism was, it seems, commentaries that were destined for teaching purposes, such as that of Ibn Qunfūdh (1330–1407) of Algeria on the algebra of Ibn al-Bannā' (1256–1321) of Morocco.

Other parts of algebra developed as well. Both Greeks and Hindus had studied indeterminate equations, and the translation of this material and the application of the newly developed algebra led to the investigation of Diophantine equations by writers like Abū Kāmil, al-Karajī, and Abū Ja'far al-Khāzin (first half of 10th century), as well as to attempts to prove a special case of what is now known as Fermat's last theorem, namely that there are no rational solutions to $x^3 + y^3 = z^3$. The great scientist Alhazen (Ibn al-Haytham [965–1041]) solved problems involving congruences by what is now called Wilson's theorem, which states that, if p is a prime, then p divides $(p-1) \times (p-2) \dots \times 2 \times 1 + 1$, and al-Baghdādī gave a variant of the idea of amicable numbers by defining two numbers to "balance" if the sums of their divisors are equal.

However, not only arithmetic and algebra but geometry too underwent extensive development. Thābit ibn Qurrah, his grandson Ibrāhīm ibn Sinān (909–946), Abū Sahl al-Kūhī (d. c. 995), and Alhazen solved problems involving the pure geometry of conic sections, including the areas and volumes of plane and solid figures formed from them, and also investigated the optical properties of mirrors made from conic sections. Ibrāhīm ibn Sinān, Abū Sahl al-Kūhī, and Alhazen used the ancient technique of analysis to reduce the solution of problems to constructions involving conic sections. (Alhazen, for example, used this method to find the point on a convex spherical mirror at which a given object is seen by a given observer.) Thābit and Ibrāhīm showed how to design the curves needed for sundials. Abū al-Wafā', whose book on the arithmetic of the scribes is mentioned above, also wrote on geometric methods needed by artisans.

In addition, in the late 10th century Abū al-Wafā' and the prince Abū Naṣr Maṣūf stated and proved theorems

Algebra

Arithmetic

Geometry

of plane and spherical geometry that could be applied by astronomers and geographers, including the laws of sines and tangents. Abū Naṣr's pupil al-Bīrūnī (973–1050), who produced a vast amount of high-quality work, was one of the masters in applying these theorems to astronomy and to such problems in mathematical geography as the determination of latitudes and longitudes, the distances between cities, and the direction from one city to another.

Omar Khayyam. The mathematician and poet Omar Khayyam was born in Neyshābūr (in modern Iran) only a few years before al-Bīrūnī's death. He later lived in Samarkand and Eṣfahān, and his brilliant work there continued many of the main lines of development in 19th-century mathematics. Not only did he discover a general method of extracting roots of arbitrary high degree, but his *Algebra* contains the first complete treatment of the solution of cubic equations. Omar did this by means of conic sections, but he declared his hope that his successors would succeed where he had failed in finding an algebraic formula for the roots.

Omar was also a part of an Islāmic tradition, which included Thābit and Alhazen, of investigating Euclid's parallel postulate. To this tradition Omar contributed the idea of a quadrilateral with two congruent sides perpendicular to the base. The parallel postulate would be proved, Omar recognized, if he could show that the remaining two angles were right angles. In this he failed, but his question about the quadrilateral became the standard way of discussing the parallel postulate.

That postulate, however, was only one of the questions on the foundations of mathematics that interested Islāmic scientists. Another was the definition of ratios. Omar Khayyam, along with others before him, felt that the theory in Euclid's Book V was logically satisfactory but intuitively unappealing, so he proved that a definition known to Aristotle was equivalent to that given in Euclid. In fact, Omar argued that ratios should be regarded as "ideal numbers," and so he conceived of a much broader system of numbers than that used since Greek antiquity, that of the positive real numbers.

Islāmic mathematics to the 15th century. In the 12th century the physician as-Samaw'al continued and completed the work of al-Karajī in algebra and also provided a systematic treatment of decimal fractions as a means of approximating irrational quantities. In his method of finding roots of pure equations, $x^n = N$, he used what is now known as Horner's method to expand the binomial $(a + y)^n$. His contemporary, Sharaf ad-Dīn aṭ-Ṭūsī, late in the 12th century provided a method of approximating the positive roots of arbitrary equations, based on an approach virtually identical to that discovered by François Viète in 16th-century France. The important step here was less the general idea than the development of the numerical algorithms necessary to effect it.

Sharaf ad-Dīn was the discoverer of a device, called the linear astrolabe, that places him in another important Islāmic mathematical tradition, one that centred on the design of new forms of the ancient astronomical instrument known as the astrolabe. The astrolabe, whose mathematical theory is based on the stereographic projection of the sphere, was invented in late antiquity, but its extensive development in Islām made it the pocket watch of the medievals. In its original form, it required a different plate of horizon coordinates for each latitude, but in the 11th century the Spanish Muslim astronomer az-Zarqallū invented a single plate that worked for all latitudes. Slightly earlier, astronomers in the East had experimented with plane projections of the sphere, and al-Bīrūnī invented such a projection that could be used to produce a map of a hemisphere. The culminating masterpiece was the astrolabe of the Syrian Ibn ash-Shāṭir (1305–75), a mathematical tool that could be used to solve all the standard problems of spherical astronomy in five different ways.

On the other hand, Muslim astronomers had developed other methods for solving these problems using the highly accurate trigonometric tables and the new trigonometric theorems they had developed. Out of these developments came the creation of trigonometry as a mathematical discipline, separate from its astronomical applications, by

Naṣir ad-Dīn aṭ-Ṭūsī at his observatory in Marāgheh in the 13th century. (It was there, too, that Naṣir ad-Dīn's pupil, Qūṭb ad-Dīn ash-Shirāzī [1236–1311], and his pupil, Kamāl ad-Dīn Fārisī, using Alhazen's great work, the *Optics*, were able to give the first mathematically satisfactory explanation of the rainbow.)

Naṣir ad-Dīn's observatory was supported by a grandson of Genghis Khan, Hülegü, who sacked Baghdad in 1258. Ulugh Beg, the grandson of the Mongol conqueror Timur, founded an observatory at Samarkand in the early years of the 15th century. Ulugh Beg was himself a good astronomer, and his five (sexagesimal) place tables of sines and tangents for every minute of arc were one of the great achievements in numerical mathematics up to his time. He was also the patron of Jamshīd al-Kāshī (d. 1429), whose work *The Reckoners' Key* summarizes most of the arithmetic of his time and includes sections on algebra and practical geometry as well. Among al-Kāshī's works are a masterful computation of the value of 2π , which, when expressed in decimal fractions, is accurate to 16 places, as well as the application of a numerical method, now known as fixed-point iteration, for solving the cubic equation with $\sin 1^\circ$ as a root. His work was indeed of a quality deserving Ulugh Beg's description as "known among the famous of the world."

Al-Kāshī lived almost five centuries after the first translations of Arabic material into Latin, and by his time the Islāmic mathematical tradition had given the West not only its first versions of many of the Greek classics but also a complete set of algorithms for Hindu-Arabic arithmetic, plane and spherical trigonometry, and the powerful tool of algebra. Although mathematical inquiry continued in Islām in the centuries after al-Kāshī's time, the mathematical centre of gravity was shifting to the West. That this was so is due, of course, in no small measure to what the Western mathematicians had learned from their Islāmic predecessors during the preceding centuries. (J.L.B.)

European mathematics during the Middle Ages and Renaissance

Until the 11th century only a small part of the Greek mathematical corpus was known in the West. Because almost no one could read Greek, what little was available came from the poor texts written in Latin in the Roman Empire, together with the very few Latin translations of Greek works. Of these the most important were the treatises by Boethius, who in about AD 500 made Latin redactions of a number of Greek scientific and logical writings. His *Arithmetic*, which was based on Nicomachus, was well known and was the means by which medieval scholars learned of Pythagorean number theory. Boethius and Cassiodorus provided the material for the part of the monastic education called the quadrivium: arithmetic, geometry, astronomy, and music theory. Together with the trivium (grammar, logic, rhetoric), these subjects formed the seven liberal arts, which were taught in the monasteries, cathedral schools, and, from the 12th century on, in the universities and which constituted the principal university instruction until modern times.

For monastic life it sufficed to know how to calculate with Roman numerals. The principal application of arithmetic was a method for determining the date of Easter, the computus, that was based on the lunar cycle of 19 solar years (*i.e.*, 235 lunar revolutions) and the 28-year solar cycle. Between the time of Bede (d. 735), when the system was fully developed, and about 1500, the computus was reduced to a series of verses that were learned by rote. Until the 12th century, geometry was largely concerned with approximate formulas for measuring areas and volumes in the tradition of the Roman surveyors. About AD 1000 the French scholar Gerbert of Aurillac, later Pope Sylvester II, introduced a type of abacus, in which numbers were represented by stones bearing Arabic numerals. Such novelties were known to very few.

The transmission of Greek and Arabic learning. In the 11th century a new phase of mathematics began with the translations from Arabic. Scholars throughout Europe went to Toledo, Córdoba, and elsewhere in Spain to translate

Euclid's
parallel
postulate

Monastic
education

Develop-
ment of
trigonom-
etry

into Latin the accumulated learning of the Muslims. Along with philosophy, astronomy, astrology, and medicine, important mathematical achievements of the Greek, Indian, and Islamic civilizations became available in the West. Particularly important were Euclid's *Elements*, the works of Archimedes, and al-Khwārizmī's treatises on arithmetic and algebra. Western texts called *algorismus* (a Latin form of the name al-Khwārizmī), introduced the Hindu-Arabic numerals and applied them in calculations. Thus modern numerals first came into use in universities and then became common among merchants and other laymen. It should be noted that, up to the 15th century, calculations were often performed with board and counters. Reckoning with Hindu-Arabic numerals was used by merchants at least from the time of Leonardo of Pisa (beginning of the 13th century) first in Italy, then in the trading cities of southern Germany and France, where *maestri d'abbaco* or *Rechenmeister* taught commercial arithmetic in the various vernaculars. Some schools were private, while others were run by the community.

The universities. Mathematics was studied from a theoretical standpoint in the universities. The universities of Paris and Oxford, which were founded relatively early (c. 1200), were centres for mathematics and philosophy. Of particular importance in these universities were the Arabic-based versions of Euclid, of which there were at least four by the 12th century. Of the numerous redactions and compendia which were made, that of Campanus (c. 1250; first printed in 1482) was easily the most popular, serving as a textbook for many generations. Such redactions of the *Elements* were made to help students not only to understand Euclid's textbook but also to handle other, particularly philosophical, questions suggested by passages in Aristotle. The ratio theory of the *Elements* provided a means of expressing the various relations of the quantities associated with moving bodies, relations that now would be expressed by formulas. Also in Euclid were to be found methods of analyzing infinity and continuity (paradoxically, because Euclid always avoided infinity).

Studies of such questions led not only to new results but also to a new approach to what is now called physics. Thomas Bradwardine, who was active in Merton College, Oxford, in the first half of the 14th century, was one of the first medieval scholars to ask whether the continuum can be divided infinitely, or whether there are smallest parts (indivisibles). Among other topics, he compared different geometric shapes in terms of the multitude of points that were assumed to compose them, and from such an approach paradoxes were generated that were not to be solved for centuries. Another fertile question stemming from Euclid concerned the angle between a circle and a line tangent to it (called the horn angle): if this angle is not zero, a contradiction quickly ensues, but, if it is zero, then, by definition, there can be no angle. For the relation of force, resistance, and the speed of the body moved by this force, Bradwardine suggested an exponential law. Nicholas Oresme (d. 1382) extended Bradwardine's ideas to fractional exponents.

Another question having to do with the quantification of qualities, the so-called latitude of forms, began to be discussed at about this time in Paris and in Merton College, Oxford. Various Aristotelian qualities (e.g., heat, density, and velocity) were assigned an intensity and extension, which were sometimes represented by the height and base (respectively) of a geometric figure. The area of the figure was then considered to represent the quantity of the quality. In the important case in which the quality is the motion of a body, the intensity its speed, and the extension its time, the area of the figure was taken to represent the distance covered by the body. Uniformly accelerated motion starting at zero velocity gives rise to a triangular figure (see Figure 11). It was proved by the Merton school that the quantity of motion in such a case is equal to the quantity of a uniform motion at the speed achieved halfway through the accelerated motion; in modern formulation: $s = \frac{1}{2}at^2$ (Merton rule). Discussions like this certainly influenced Galileo indirectly and may have influenced the founding of coordinate geometry in the 17th century. Another important development in the scholas-



Figure 11: Uniformly accelerated motion; s = speed, a = acceleration, t = time, and v = velocity.

tic "calculations" was the summation of infinite series.

Basing his work on translated Greek sources, the German mathematician and astronomer Regiomontanus wrote the first book in the West on plane and spherical trigonometry independent of astronomy in about 1464 (printed in 1533). He also published tables of sines and tangents that were in constant use for more than two centuries.

The Renaissance. Italian artists and merchants influenced the mathematics of the late Middle Ages and the Renaissance in several ways. In the 15th century a group of Tuscan artists, including Filippo Brunelleschi, Leon Battista Alberti, and Leonardo da Vinci, incorporated linear perspective into their practice and teaching, about a century before the subject was formally treated by mathematicians. Italian *maestri d'abbaco* tried, albeit unsuccessfully, to solve nontrivial cubic equations. In fact, the first general solution was found by Scipione Del Ferro at the beginning of the 16th century and rediscovered by Niccolò Tartaglia several years later. The solution was published by Girolamo Cardano in his *Ars magna* in 1545, together with Lodovico Ferrari's solution of the quartic equation.

By 1380 an algebraic symbolism had been developed in Italy in which letters were used for the unknown, for its square, and for constants. The symbols used today for the unknown (for example, x), the square root sign, and the signs $+$ and $-$ came into general use in southern Germany beginning in about 1450. They were used by Regiomontanus and by Fridericus Gerhart and received an impetus in about 1486 at the University of Leipzig from Johann Widman. The idea of distinguishing between known and unknown quantities in algebra was first consistently applied by François Viète, with vowels for unknown and consonants for known quantities. Viète found some relations between the coefficients of an equation and its roots. This was suggestive of the idea, explicitly stated by Albert Girard in 1629 and proved by Gauss in 1799, that an equation of degree n has n roots. Complex numbers, which are implicit in such ideas, were gradually accepted about the time of Rafael Bombelli (d. 1572), who used them in connection with the cubic.

Apollonius' *Conics* and the investigations of areas (quadratures) and of volumes (cubatures) of Archimedes formed part of the humanistic learning of the 16th century. These studies strongly influenced the later developments of analytic geometry, the infinitesimal calculus, and the theory of functions, subjects that were developed in the 17th century. (Me.F.)

Mathematics in the 17th and 18th centuries

THE 17TH CENTURY

The 17th century, the period of the scientific revolution, witnessed the consolidation of Copernican heliocentric astronomy and the establishment of inertial physics in the work of Kepler, Galileo, Descartes, and Newton. This period was also one of intense activity and innovation in mathematics. Advances in numerical calculation, the development of symbolic algebra and analytic geometry, and the invention of the differential and integral calculus resulted in a major expansion of the subject areas of mathematics. By the end of the 17th century a program of research based in analysis had replaced classical Greek geometry at the centre of advanced mathematics. In the next century this program would continue to develop in close association with physics, more particularly mechanics and theoretical astronomy. The extensive use of

Influence
of Euclid's
Elements

The
Merton
rule

Develop-
ment of
analysis

analytic methods, the incorporation of applied subjects, and the adoption of a pragmatic attitude to questions of logical rigour distinguished the new mathematics from traditional geometry.

Institutional background. Until the middle of the 17th century, mathematicians worked alone or in small groups, publishing their work in books or communicating with other researchers by letter. At a time when people were often slow to publish, “invisible colleges,” networks of scientists who corresponded privately, played an important role in coordinating and stimulating mathematical research. Marin Mersenne in Paris acted as a clearinghouse for new results, informing his many correspondents—including Fermat, Descartes, Blaise Pascal, Gilles Personne de Roberval, and Galileo—of challenge problems and novel solutions. Later in the century John Collins, librarian of London’s Royal Society, performed a similar function among British mathematicians.

In 1660 the Royal Society of London was founded, to be followed in 1666 by the Academy of Sciences in France, in 1700 by the Berlin Academy, and in 1724 by the St. Petersburg Academy. The official publications sponsored by the academies, as well as independent journals such as the *Acta Eruditorum* (founded in 1682), made possible the open and prompt communication of research findings. Although universities in the 17th century provided some support for mathematics, they became increasingly ineffective as state-supported academies assumed direction of advanced research.

Numerical calculation. The development of new methods of numerical calculation was a response to the increased practical demands of numerical computation, particularly in trigonometry, navigation, and astronomy. New ideas spread quickly across Europe and resulted by 1630 in a major revolution in numerical practice.

Simon Stevin of Holland, in his short pamphlet *La Disme* (1585), introduced decimal fractions to Europe and showed how to extend the principles of Hindu-Arabic arithmetic to calculation with these numbers. Stevin emphasized the utility of decimal arithmetic “for all accounts that are encountered in the affairs of men,” and he explained in an appendix how it could be applied to surveying, stereometry, astronomy, and mensuration. His idea was to extend the base-10 positional principle to numbers with fractional parts, with a corresponding extension of notation to cover these cases. In his system the number 237.578 was denoted

237 ④ 5 ① 7 ② 8 ③,

in which the digits to the left of the zero are the integral part of the number. To the right of the zero are the digits of the fractional part, with each digit succeeded by a circled number that indicates the negative power to which 10 is raised. Stevin showed how the usual arithmetic of whole numbers could be extended to decimal fractions using rules that determined the positioning of the negative powers of 10.

In addition to its practical utility, *La Disme* was significant for the way it undermined the dominant style of classical Greek geometry in theoretical mathematics. Stevin’s proposal required a rejection of the distinction in Euclidean geometry between magnitude, which is continuous, and number, which is a multitude of indivisible units. For Euclid, unity, or one, was a special sort of thing, not number but the origin, or principle, of number. The introduction of decimal fractions seemed to imply that the unit could be subdivided and that arbitrary continuous magnitude could be represented numerically; it implicitly supposed the concept of a general positive real number.

Tables of logarithms were first published in 1614 by the Scottish baron John Napier in his treatise *Mirifici Logarithmorum Canonis Descriptio* (*Description of the Marvelous Canon of Logarithms*). This work was followed (posthumously) five years later by another in which Napier set forth the principles used in the construction of his tables. The basic idea behind logarithms is that addition and subtraction are easier to perform than multiplication and division, which, as Napier observed, require a “te-

dious expenditure of time” and are subject to “slippery errors.” By the law of exponents, $a^n a^m = a^{n+m}$, that is, in the multiplication of numbers the exponents are related additively. By correlating the geometric sequence of numbers a, a^2, a^3, \dots (a is called the base) and the arithmetic sequence $1, 2, 3, \dots$ and interpolating to fractional values, it is possible to reduce the problem of multiplication and division to one of addition and subtraction. To do this Napier chose a base that was very close to 1, differing from it by only $1/10^7$. The resulting geometric sequence therefore yielded a dense set of values, suitable for constructing a table.

In his work of 1619 Napier presented an interesting kinematic model to generate the geometric and arithmetic sequences used in the construction of his tables. Assume two particles move along separate lines from given initial points. The particles begin moving at the same instant with the same velocity. The first particle continues to move with a speed that is decreasing, proportional at each instant to the distance remaining between it and some given fixed point on the line. The second particle moves with a constant speed equal to its initial velocity. Given any increment of time, the distances traveled by the first particle in successive increments form a geometrically decreasing sequence. The corresponding distances traveled by the second particle form an arithmetically increasing sequence. Napier was able to use this model to derive theorems yielding precise limits to approximate values in the two sequences.

Napier’s kinematic model indicated how skilled mathematicians had become by the early 17th century in analyzing nonuniform motion. Kinematic ideas, which appeared frequently in mathematics of the period, provided a clear and visualizable means for the generation of geometric magnitude. The conception of a curve traced by a particle moving through space later played a significant role in the development of the calculus.

Napier’s ideas were taken up and revised by the English mathematician Henry Briggs, the first Savilian professor of geometry at Oxford. In 1624 Briggs published an extensive table of common logarithms, or logarithms to the base 10. Because the base was no longer close to 1, the table could not be obtained as simply as Napier’s, and Briggs therefore devised techniques involving the calculus of finite differences to facilitate calculation of the entries. He also devised interpolation procedures of great computational efficiency to obtain intermediate values.

In Switzerland the instrument maker Joost Bürgi arrived at the idea for logarithms independently of Napier, although he did not publish his results until 1620. Four years later a table of logarithms prepared by Johannes Kepler appeared in Marburg. Both Bürgi and Kepler were astronomical observers, and Kepler included logarithmic tables in his famous Rudolphine Tables, astronomical tabulations of planetary motion derived using the assumption of elliptical orbits about the Sun.

Analytic geometry. The invention of analytic geometry was, next to the differential and integral calculus, the most important mathematical development of the 17th century. Originating in the work of the French mathematicians Viète, Fermat, and Descartes, it had by the middle of the century established itself as a major program of mathematical research.

Two tendencies in contemporary mathematics stimulated the rise of analytic geometry. The first was an increased interest in curves, resulting in part from the recovery and Latin translation of the classical treatises of Apollonius, Archimedes, and Pappus, and in part from the increasing importance of curves in such applied fields as astronomy, mechanics, optics, and stereometry. The second was the emergence a century earlier of an established algebraic practice in the work of the Italian and German algebraists and its subsequent shaping into a powerful mathematical tool by Viète at the end of the century.

Viète was a prominent representative of the humanist movement in mathematics that set itself the project of restoring and furthering the achievements of the classical Greek geometers. In his *In artem analyticam isagoge* (“Introduction to the Analytic Arts”; 1591) Viète, as part

Influence of Greek geometers

of his program of rediscovering the method of analysis used by the ancient Greek mathematicians, proposed new algebraic methods that employed variables, constants, and equations, but he saw this as an advancement over the ancient method, a view he arrived at by comparing the geometric analysis contained in Book VII of Pappus' *Collection* with the arithmetic analysis of Diophantus' *Arithmetica*. Pappus had employed an analytic method for the discovery of theorems and the construction of problems; in analysis, by contrast to synthesis, one proceeds from what is sought until one arrives at something known. In approaching an arithmetic problem by laying down an equation among known and unknown magnitudes and then solving for the unknown, one was, Viète reasoned, following an "analytic" procedure.

Viète's notation

Viète introduced the concept of algebraic variable, which he denoted using a capital vowel (*A, E, I, O, U*) as well as the concept of parameter (an unspecified constant quantity), denoted by a capital consonant (*B, C, D*, and so on). In his system the equation $5BA^2 - 2CA + A^3 = D$ would appear as *B5* in *A* quad $- C$ plano 2 in *A* + *A* cub aequatur *D* solido.

Viète retained the classical principle of homogeneity, according to which terms added together must all be of the same dimension. In the above equation, for example, each of the terms has the dimension of a solid or cube; thus the constant *C*, which denotes a plane, is combined with *A* to form a quantity having the dimension of a solid.

It should be noted that in Viète's scheme the symbol *A* is part of the expression for the object obtained by operating on the magnitude denoted by *A*. Thus operations on the quantities denoted by the variables are reflected in the algebraic notation itself. This innovation, considered by historians of mathematics to be a major conceptual advance in algebra, facilitated the study of the symbolic solution of algebraic equations and led to the creation of the first conscious theory of equations.

After Viète's death the analytic art was applied to the study of curves by his countrymen Fermat and Descartes. Both men were motivated by the same goal, to apply the new algebraic techniques to Apollonius' theory of loci as preserved in Pappus' *Collection*. The most celebrated of these problems consisted of finding the curve or locus traced by a point whose distances from several fixed lines satisfied a given relation.

Fermat adopted Viète's notation in his paper of 1636, "Ad Locos Planos et Solidos Isagoge" ("Introduction to Plane and Solid Loci"). The title of the paper refers to the ancient classification of curves as plane (straight lines, circles), solid (ellipses, parabolas, and hyperbolas), or linear (curves defined kinematically or by a locus condition). Fermat considered an equation among two variables. One of the variables represented a line measured horizontally from a given initial point, while the other represented a second line positioned at the end of the first line and inclined at a fixed angle to the horizontal. As the first variable varied in magnitude, the second took on a value determined by the equation, and the endpoint of the second line traced out a curve in space. By means of this construction Fermat was able to formulate the fundamental principle of analytic geometry:

Whenever two unknown quantities are found in final equality, there results a locus fixed in place, and the endpoint of one of these unknown quantities describes a straight line or a curve.

The principle implied a correspondence between two different classes of mathematical objects, geometric curves and algebraic equations. In the paper of 1636 Fermat showed that, if the equation is a quadratic, then the curve is a conic section, that is, an ellipse, parabola, or hyperbola. He also showed that the determination of the curve given by an equation is simplified by a transformation involving a change of variables to an equation in standard form.

Descartes's *La Géométrie* appeared in 1637 as an appendix to his famous *Discours de la méthode*, the treatise that presented the foundation of his philosophical system. Although supposedly an example from mathematics of his rational method, *La Géométrie* was a technical treatise understandable independently of philosophy. It was des-

igned to become one of the most influential books in the history of mathematics.

In the opening sections of *La Géométrie* Descartes introduced two innovations. In place of Viète's notation he initiated the modern practice of denoting variables by letters at the end of the alphabet (*x, y, z*) and parameters by letters at the beginning of the alphabet (*a, b, c*) and of using exponential notation to indicate powers of *x* (x^2, x^3, \dots). More significantly conceptually, he set aside Viète's principle of homogeneity, showing by means of a simple construction how to represent multiplication and division of lines by lines; thus all magnitudes (lines, areas, and volumes) could be represented independently of their dimension in the same way.

Descartes's goal in *La Géométrie* was to achieve the construction of solutions to geometric problems by means of instruments that were acceptable generalizations of ruler and compass. Algebra was a tool to be used in this program:

If, then, we wish to solve any problem, we first suppose the solution already effected, and give names to all the lines that seem necessary for its construction—to those that are unknown as well as to those that are known. Then, making no distinction in any way between known and unknown lines, we must unravel the difficulty in any way that shows most naturally the relations between these lines, until we find it possible to express a single quantity in two ways. This will constitute an equation; since the terms of one of these two expressions are together equal to the terms of the other.

In the problem of Apollonius, for example, one sought to find the locus of points whose distances from a collection of fixed lines satisfied a given relation. One used this relation to derive an equation and then obtained points on the curve given by the roots of the equation using a geometric procedure involving acceptable instruments of construction.

Descartes described instruments more general than the compass for drawing "geometric" curves. He stipulated that the parts of the instrument be linked together so that the ratio of the motions of the parts could be knowable. This restriction excluded "mechanical" curves generated by kinematic processes. The Archimedean spiral, for example, was generated by a point moving on a line as the line rotated uniformly about the origin. The ratio of the circumference to the diameter did not permit exact determination:

the ratios between straight and curved lines are not known, and I even believe cannot be discovered by men, and therefore no conclusion based upon such ratios can be accepted as rigorous and exact.

Descartes concluded that a geometric or nonmechanical curve was one whose equation $f(x, y) = 0$ was a polynomial of finite degree in two variables. He wished to restrict mathematics to the consideration of such curves.

Descartes's emphasis on construction reflected his classical orientation. His conservatism with respect to what curves were acceptable in mathematics further distinguished him as a traditional thinker. At the time of his death in 1650, he had been overtaken by events, as research moved away from questions of construction to problems of finding areas (then called problems of quadrature) and tangents. The geometric objects that were then of growing interest were precisely the mechanical curves that Descartes had wished to banish from mathematics.

Following the important results achieved in the 16th century by Cardano and the Italian algebraists, the theory of algebraic equations reached an impasse. The ideas needed to investigate equations of degree higher than four were slow to develop. The immediate historical influence of Viète, Fermat, and Descartes was to furnish algebraic methods for the investigation of curves. A vigorous school of research became established in Leiden around Frans van Schooten, a Dutch mathematician who edited and published in 1649 a Latin translation of *La Géométrie*. Van Schooten published a second two-volume translation of the same work in 1659–1661 that also contained mathematical appendices by three of his disciples, Johan de Witt, Johan Hudde, and Hendrick van Heuraet. The Leiden group of mathematicians, which also included

La Géométrie

Geometric curves

Christiaan Huygens, was in large part responsible for the rapid development of Cartesian geometry in the middle of the century.

The calculus. The historian Carl Boyer has called the calculus "the most effective instrument for scientific investigation that mathematics has ever produced." As the mathematics of variability and change, the calculus was the characteristic product of the scientific revolution. The subject was properly the invention of two mathematicians, the German Gottfried Wilhelm Leibniz and the Englishman Isaac Newton. Both men published their researches in the 1680s, Leibniz in 1684 in the recently founded journal *Acta Eruditorum* and Newton in 1687 in his great treatise *Principia Mathematica*. Although a bitter dispute over priority developed later among followers of the two men, it is now clear that they each arrived at the calculus independently.

The calculus developed from techniques to solve two types of problems, the determination of areas and volumes and the calculation of tangents to curves. In classical geometry Archimedes had advanced furthest in this part of mathematics, having used the method of exhaustion to establish rigorously various results on areas and volumes and having derived for some curves (e.g., the spiral) significant results concerning tangents. In the early 17th century there was a sharp revival of interest in both classes of problems. The decades between 1610 and 1670, referred to in the history of mathematics as "the precalculus period," were at a time of remarkable activity in which researchers throughout Europe contributed novel solutions and competed with each other to arrive at important new methods.

The precalculus period. In his treatise *Geometria Indivisibilibus Continuatorum* (1635) (*Geometry by Indivisibles of Continuum*) Cavalieri, a professor of mathematics at the University of Bologna, formulated a systematic method for the determination of areas and volumes. As had Archimedes, Cavalieri regarded a plane figure as being composed of a collection of indivisible lines, "all the lines" of the plane figure. The collection was generated by a fixed line moving through space parallel to itself. Cavalieri showed that these collections could be interpreted as magnitudes obeying the rules of Euclidean ratio theory. In proposition 4 of Book II, he derived the result that is written today as

$$\int_0^1 x^2 dx = \frac{1}{3} :$$

Let there be given a parallelogram in which a diagonal is drawn; then "all the squares" of the parallelogram will be triple "all the squares" of each of the triangles determined by the diagonal.

Cavalieri showed that this proposition could be interpreted in different ways, as asserting that the volume of a cone is one-third the volume of the circumscribed cylinder or that the area under a segment of a parabola is one-third the area of the associated rectangle. In a later treatise he generalized the result by proving

$$\int_0^1 x^n dx = \frac{1}{(n+1)}$$

for $n=3$ to $n=9$. To establish these results he introduced transformations among the variables of the problem, using a result equivalent to the binomial theorem for integral exponents. The ideas involved went beyond anything that had appeared in the classical Archimedean theory of content.

Although Cavalieri was successful in formulating a systematic method based on general concepts, his ideas were not easy to apply. The derivation of very simple results required intricate geometric considerations, and the turgid style of the *Geometria Indivisibilibus* was a barrier to its reception.

John Wallis presented a quite different approach to the theory of quadratures in his *Arithmetica Infinitorum* (1655) (*Arithmetic of Infinites*). Wallis, a successor to Henry Briggs as the Savilian professor of geometry at Oxford, was a champion of the new methods of arithmetic algebra that he had learned from his teacher William

Oughtred. Wallis expressed the area under a curve as the sum of an infinite series and used clever and unrigorous inductions to determine its value. To calculate the area under the parabola,

$$\int_0^1 x^2 dx,$$

he considered the successive sums

$$\frac{0+1}{1+1} = \frac{1}{3} + \frac{1}{6}, \frac{0+1+4}{4+4+4} = \frac{1}{3} + \frac{1}{12}, \frac{0+1+4+9}{9+9+9+9} = \frac{1}{3} + \frac{1}{18}$$

and inferred by "induction" the general relation

$$\frac{0^2 + 1^2 + 2^2 \dots + n^2}{n^2 + n^2 + n^2 \dots + n^2} = \frac{1}{3} + \frac{1}{6n}.$$

By letting the number of terms be infinite, he obtained $\frac{1}{3}$ as the limiting value of the expression. With more complicated curves he achieved very impressive results, including the infinite expression now known as Wallis' product:

$$\frac{4}{\pi} = \frac{3}{2} \cdot \frac{3}{4} \cdot \frac{5}{4} \cdot \frac{5}{6} \cdot \frac{7}{6} \dots$$

Research on the determination of tangents, the other subject leading to the calculus, proceeded along different lines. In *La Géométrie* Descartes had presented a method that could in principle be applied to any algebraic or "geometric" curve—i.e., any curve whose equation was a polynomial of finite degree in two variables. The method depended upon finding the normal, the line perpendicular to the tangent, using the algebraic condition that it be the unique radius to intersect the curve in only one point. Descartes's method was simplified by Hudde, a member of the Leiden group of mathematicians, and was published in 1659 in van Schooten's edition of *La Géométrie*.

A class of curves of growing interest in the 17th century were those generated kinematically by a point moving through space. The famous cycloidal curve, for example, was traced by a point on the perimeter of a wheel that rolled on a line without slipping or sliding. These curves were nonalgebraic and hence could not be treated by Descartes's method. Gilles de Roberval, professor at the Collège Royale in Paris, devised a method borrowed from dynamics to determine their tangents. In his analysis of projectile motion Galileo had shown that the instantaneous velocity of a particle is compounded of two separate motions: a constant horizontal motion and an increasing vertical motion due to gravity. If the motion of the generating point of a kinematic curve is likewise regarded as the sum of two velocities, then the tangent will lie in the direction of their sum. Roberval applied this idea to several different kinematic curves, obtaining results that were often ingenious and elegant.

Non-algebraic curves

In an essay of 1636 circulated among French mathematicians, Fermat presented a method of tangents adapted from a procedure he had devised to determine maxima and minima and used it to find tangents to several algebraic curves of the form $y = x^n$. His account was short and contained no explanation of the mathematical basis of the new method. It is possible to see in his procedure an argument involving infinitesimals, and Fermat has sometimes been proclaimed as the discoverer of the differential calculus. Modern historical study, however, suggests that he was working with concepts introduced by Viète and that his method was based on finite algebraic ideas.

Isaac Barrow, the Lucasian professor of mathematics at Cambridge, published in 1670 his *Lectiones Geometricae* (*Geometrical Lectures*), a treatise that more than any other anticipated the unifying ideas of the calculus. In it he adopted a purely geometric form of exposition to show how the determinations of areas and tangents are inverse problems. He began with a curve and considered the slope of its tangent corresponding to each value of the abscissa. He then defined an auxiliary curve by the condition that its ordinate be equal to this slope and showed that the area under the auxiliary curve corresponding to a given

Cavalieri

abscissa is equal to the rectangle whose sides are unity and the ordinate of the original curve. When reformulated analytically, this result expresses the inverse character of differentiation and integration, the fundamental theorem of the calculus. Although Barrow's decision to proceed geometrically prevented him from taking the final step to a true calculus, his lectures influenced both Newton and Leibniz.

Newton and Leibniz. The essential insight of Newton and Leibniz was to use Cartesian algebra to synthesize the earlier results and to develop algorithms that could be applied uniformly to a wide class of problems. The formative period of Newton's researches was from 1665 to 1670, while Leibniz worked a few years later, in the 1670s. Their contributions differ in origin, development, and influence, and it is necessary to consider each man separately.

The son of an English farmer, Newton became in 1669 the Lucasian professor of mathematics at Cambridge University. Newton's earliest researches in mathematics grew in 1665 from his study of van Schooten's edition of *La Géométrie* and Wallis' *Arithmetica Infinitorum*. Using the Cartesian equation of the curve, he reformulated Wallis' results, introducing for this purpose infinite sums in the powers of an unknown x , now known as infinite series. Possibly under the influence of Barrow, he used infinitesimals to establish for various curves the inverse relationship of tangents and areas. The operations of differentiation and integration emerged in his work as analytic processes that could be applied generally to investigate curves.

Unusually sensitive to questions of rigour, Newton at a fairly early stage tried to establish his new method on a sound foundation using ideas from kinematics. A variable was regarded as a "fluent," a magnitude that flows with time; its derivative or rate of change with respect to time was called a "fluxion," denoted by the given variable with a dot above it. The basic problem of the calculus was to investigate relations among fluents and their fluxions. Newton finished a treatise on the method of fluxions as early as 1671, although it was not published until 1736. In the 18th century this method became the preferred approach to the calculus among British mathematicians, especially after the appearance in 1742 of Colin Maclaurin's influential *Treatise of Fluxions*.

Newton first published the calculus in Book I of his great *Philosophiæ Naturalis Principia Mathematica* (1687) (*Mathematical Principles of Natural Philosophy*). Originating as a treatise on the dynamics of particles, the *Principia* presented an inertial physics that combined Galileo's mechanics and Kepler's planetary astronomy. It was written in the early 1680s at a time when Newton was reacting against Descartes's science and mathematics. Setting aside the analytic method of fluxions, Newton introduced in 11 introductory lemmas his calculus of first and last ratios, a geometric theory of limits that provided the mathematical basis of his dynamics.

Newton's use of the calculus in the *Principia* is illustrated by proposition 11 of Book I: If the orbit of a particle moving under a centripetal force is an ellipse with the centre of force at one focus, then the force is inversely proportional to the square of the distance from the centre. Because the planets were known by Kepler's laws to move in ellipses with the Sun at one focus, this result supported his inverse square law of gravitation. To establish the proposition, Newton derived an approximate measure for the force by using small lines defined in terms of the radius (the line from the force centre to the particle) and the tangent to the curve at a point. This result expressed geometrically the proportionality of force to vector acceleration. Using properties of the ellipse known from classical geometry, Newton calculated the limit of this measure and showed that it was equal to a constant times one over the square of the radius.

Newton avoided analytical processes in the *Principia* by expressing magnitudes and ratios directly in terms of geometric quantities, both finite and infinitesimal. His decision to eschew analysis constituted a striking rejection of the algebraic methods that had been important in his own early researches on the calculus. Although the *Principia* was of inestimable value for later mechanics, it would be

reworked by researchers on the Continent and expressed in the mathematical idiom of the Leibnizian calculus.

Leibniz's interest in mathematics was aroused in 1672 during a visit to Paris, where the Dutch mathematician Christiaan Huygens introduced him to his work on the theory of curves. Under Huygens' tutelage Leibniz immersed himself for the next several years in the study of mathematics. He investigated relationships among the summing and differencing of finite and infinite sequences of numbers. Having read Barrow's geometric lectures, he devised a transformation rule to calculate quadratures, obtaining the famous infinite series for $\pi/4$:

$$\frac{\pi}{4} = \frac{1}{1} - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \dots$$

Leibniz was interested in questions of logic and notation, of how to construct a *characteristica universalis* for rational investigation. After considerable experimentation he arrived by the late 1670s at an algorithm based on the symbols d and \int . He first published his research on differential calculus in 1684 in an article in the *Acta Eruditorum*, "Nova Methodus pro Maximis et Minimis, Itemque Tangentibus, qua nec Fractas nec Irrationales Quantitates Moratur, et Singulare pro illi Calculi Genus" ("A New Method for Maxima and Minima as Well as Tangents, Which Is Impeded Neither by Fractional nor by Irrational Quantities, and a Remarkable Type of Calculus for This"). In this article he introduced the differential dx satisfying the rules $d(x+y) = dx + dy$ and $d(xy) = xdy + ydx$ and illustrated his calculus with a few examples. Two years later he published a second article, "On a Deeply Hidden Geometry," in which he introduced and explained the symbol \int for integration. He stressed the power of his calculus to investigate transcendental curves, the very class of "mechanical" objects Descartes had believed lay beyond the power of analysis, and derived a simple analytical formula for the cycloid.

Leibniz continued to publish results on the new calculus in the *Acta Eruditorum* and began to explore his ideas in extensive correspondence with other scholars. Within a few years he had attracted a group of researchers to promulgate his methods, including the brothers Johann and Jakob Bernoulli in Basel and the priest Pierre Varignon and Guillaume-François-Antoine de L'Hospital in Paris. In 1700 he convinced Frederick III of Prussia to establish the Brandenburg Society of Sciences (later renamed the Berlin Academy of Sciences), with himself appointed president for life.

Leibniz's vigorous espousal of the new calculus, the didactic spirit of his writings, and his ability to attract a community of researchers contributed to his enormous influence on subsequent mathematics. In contrast, Newton's slowness to publish and his personal reticence resulted in a reduced presence within European mathematics. Although the British school in the 18th century included capable researchers, Abraham De Moivre, James Stirling, Brook Taylor, and Maclaurin among them, they failed to establish a program of research comparable to that established by Leibniz's followers on the Continent. There is a certain tragedy in Newton's isolation and his reluctance to acknowledge the superiority of continental analysis. As the historian Michael Mahoney observed:

Whatever the revolutionary influence of the *Principia*, mathematics would have looked much the same if Newton had never existed. In that endeavour he belonged to a community, and he was far from indispensable to it.

THE 18TH CENTURY

Institutional background. After 1700 a movement to found learned societies on the model of Paris and London spread throughout Europe and the American colonies. The academy was the predominant institution of science until it was displaced by the university in the 19th century. The leading mathematicians of the period, such as Leonhard Euler, Jean Le Rond d'Alembert, and Joseph-Louis Lagrange, pursued academic careers at St. Petersburg, Paris, and London.

The Paris Academy of Sciences provides an informative study of the 18th-century learned society. The academy was divided into six sections, three for the mathematical and three for the physical sciences. The mathematical sections were for geometry, astronomy, and mechanics, the physical sections for chemistry, anatomy, and botany. Membership in the academy was divided by section, with each section contributing three pensionnaires, two associates, and two adjuncts. There was also a group of free associates, distinguished men of science from the provinces, and foreign associates, eminent international figures in the field. A larger group of 70 corresponding members had partial privileges, including the right to communicate reports to the academy. The administrative core consisted of a permanent secretary, treasurer, president, and vice president. In a given year the average total membership in the academy was 153.

Prominent characteristics of the academy included its small and elite membership, made up heavily of men from the middle class, and its emphasis on the mathematical sciences. In addition to holding regular meetings and publishing memoirs, the academy organized scientific expeditions and administered prize competitions on important mathematical and scientific questions.

The historian Roger Hahn has noted that the academy in the 18th century allowed "the coupling of relative doctrinal freedom on scientific questions with rigorous evaluations by peers," an important characteristic of modern professional science. Academic mathematics and science did, however, foster a stronger individualistic ethos than is usual today. A determined individual such as Euler or Lagrange could emphasize a given program of research through his own work, the publications of the academy, and the setting of the prize competitions. The academy as an institution may have been more conducive to the solitary patterns of research in a theoretical subject like mathematics than it was to the experimental sciences. The separation of research from teaching is perhaps the most striking characteristic that distinguished the academy from the model of university-based science which developed in the 19th century.

Analysis and mechanics. The scientific revolution had bequeathed to mathematics a major program of research in analysis and mechanics. The period from 1700 to 1800, "the century of analysis," witnessed the consolidation of the calculus and its extensive application to mechanics. With expansion came specialization, as different parts of the subject acquired their own identity: ordinary and partial differential equations, calculus of variations, infinite series, and differential geometry. The applications of analysis were also varied, including the theory of the vibrating string, particle dynamics, the theory of rigid bodies, the mechanics of flexible and elastic media, and the theory of compressible and incompressible fluids. Analysis and mechanics developed in close association, with problems in one giving rise to concepts and techniques in the other, and all the leading mathematicians of the period made important contributions to mechanics.

The close relationship between mathematics and mechanics in the 18th century had roots extending deep into Enlightenment thought. In the organizational chart of knowledge at the beginning of the preliminary discourse to the *Encyclopédie*, d'Alembert distinguished between "pure" mathematics (geometry, arithmetic, algebra, calculus) and "mixed" mathematics (mechanics, geometric astronomy, optics, art of conjecturing). Mathematics generally was classified as a "science of nature" and separated from logic, a "science of man." The modern disciplinary division between physics and mathematics and the association of the latter with logic had not yet developed.

Mathematical mechanics itself as it was practiced in the 18th century differed in important respects from later physics. The goal of modern physics is to explore the ultimate particulate structure of matter and to arrive at fundamental laws of nature to explain physical phenomena. The character of applied investigation in the 18th century was rather different. The material parts of a given system and their interrelationship were idealized for the purposes of analysis. A material object could be treated

as a point-mass (a mathematical point at which it is assumed all the mass of the object is concentrated), as a rigid body, as a continuously deformable medium, and so on. The intent was to obtain a mathematical description of the macroscopic behaviour of the system rather than to ascertain the ultimate physical basis of the phenomena. In this respect the 18th-century viewpoint is closer to modern mathematical engineering than it is to physics.

Mathematical research in the 18th century was coordinated by the Paris, Berlin, and St. Petersburg academies, as well as by several smaller provincial scientific academies and societies. Although England and Scotland were important centres early in the century, with Maclaurin's death in 1746 the British flame was all but extinguished.

History of analysis. The history of analysis in the 18th century can be followed in the official memoirs of the academies and in independently published expository treatises. In the first decades of the century the calculus was cultivated in an atmosphere of intellectual excitement, as mathematicians applied the new methods to a range of problems in the geometry of curves. The brothers Johann and Jakob Bernoulli showed that the shape of a smooth wire along which a particle descends in the least time is the cycloid, a transcendental curve much studied in the previous century. Working in a spirit of keen rivalry, the two brothers arrived at ideas that would later develop into the calculus of variations. In his study of the rectification of the lemniscate, a ribbon-shaped curve discovered by Jakob Bernoulli in 1694, Giulio Carlo Fagnano (1682–1766) introduced ingenious analytic transformations that laid the foundation for the theory of elliptic integrals. Nikolaus I Bernoulli (1687–1759), the nephew of Johann and Jakob, proved the equality of mixed second-order partial derivatives and made important contributions to differential equations by the construction of orthogonal trajectories to families of curves. Pierre Varignon (1654–1722), Johann Bernoulli, and Jakob Hermann (1678–1733) continued to develop analytic dynamics as they adapted Leibniz's calculus to the inertial mechanics of Newton's *Principia*.

Geometric conceptions and problems predominated in the early calculus. This emphasis on the curve as the object of study provided coherence to what was otherwise a disparate collection of analytic techniques. With its continued development, the calculus gradually became removed from its origins in the geometry of curves, and a movement emerged to establish the subject on a purely analytic basis. In a series of textbooks published in the middle of the century, the Swiss mathematician Leonhard Euler systematically accomplished the separation of the calculus from geometry. In his *Introductio in Analysis In-finitorum* (1748; *Introduction to the Analysis of Infinities*) he made the notion of function the central organizing concept of analysis:

A function of a variable quantity is an analytical expression composed in any way from the variable and from numbers or constant quantities.

Euler's analytic approach is illustrated by his introduction of the sine and cosine functions. Trigonometric tables had existed since antiquity, and the relations between sines and cosines were commonly used in mathematical astronomy. In the early calculus mathematicians had derived in their study of periodic mechanical phenomena the differential equation

$$\frac{dy}{dx} = \frac{1}{\sqrt{1-x^2}}$$

and they were able to interpret its solution geometrically in terms of lines and angles in the circle. Euler was the first to introduce the sine and cosine functions as quantities whose relation to other quantities could be studied independently of any geometric diagram.

Euler's analytic approach to the calculus received support from his younger contemporary Joseph-Louis Lagrange, who, following Euler's death in 1783, replaced him as the leader of European mathematics. In 1755 the 19-year-old Lagrange wrote to Euler to announce the discovery of a new algorithm in the calculus of variations, a subject to

which Euler had devoted an important treatise 11 years earlier. Euler had used geometric ideas extensively and had emphasized the need for analytic methods. Lagrange's idea was to introduce the new symbol δ into the calculus and to experiment formally until he had devised an algorithm to obtain the variational equations. Mathematically quite distinct from Euler's procedure, his method required no reference to the geometric configuration. Euler immediately adopted Lagrange's idea, and in the next several years the two men systematically revised the subject using the new techniques.

In 1766 Lagrange was invited by the Prussian king, Frederick II the Great, to become mathematics director of the Berlin Academy. During the next two decades he wrote important memoirs on nearly all of the major areas of mathematics. In 1788 he published his famous *Mécanique analytique*, a treatise that used variational ideas to present mechanics from a unified analytic viewpoint. In the preface Lagrange wrote:

One will find no Figures in this Work. The methods that I present require neither constructions nor geometrical or mechanical reasonings, but only algebraic operations, subject to a regular and uniform course. Those who admire Analysis, will with pleasure see Mechanics become a new branch of it, and will be grateful to me for having extended its domain.

Following the death of Frederick the Great, Lagrange traveled to Paris to become a pensionnaire of the Academy of Sciences. With the establishment of the École Polytechnique (French: "Polytechnic School") in 1794 he was asked to deliver the lectures on mathematics. There was a concern in European mathematics at the time to place the calculus on a sound basis, and Lagrange used the occasion to develop his ideas for an algebraic foundation of the subject. The lectures were published in 1797 under the title *Théorie des fonctions analytiques* ("Theory of Analytical Functions"), a treatise whose contents were summarized in its longer title "Containing the Principles of the Differential Calculus Disengaged from All Consideration of Infinitesimals, Vanishing Limits or Fluxions and Reduced to the Algebraic Analysis of Finite Quantities." Lagrange published a second treatise on the subject in 1801, a work that appeared in a revised and expanded form in 1806.

The range of subjects presented and the consistency of style distinguished Lagrange's didactic writings from other contemporary expositions of the calculus. He began with Euler's notion of a function as an analytic expression composed of variables and constants. He defined the "derived function" or derivative $f'(x)$ of $f(x)$ to be the coefficient of i in the Taylor expansion of $f(x+i)$. Assuming the general possibility of such expansions, he attempted a rather complete theory of the differential and integral calculus, including extensive applications to geometry and mechanics. Lagrange's lectures represented the most advanced development of the 18th-century analytic conception of the calculus.

Beginning with Augustin-Louis Cauchy in the 1820s, later mathematicians used the concept of limit to establish the calculus on an arithmetic basis. The algebraic viewpoint of Euler and Lagrange was rejected. To arrive at a proper historical appreciation of their work it is necessary to reflect on the meaning of analysis in the 18th century. Since Viète, analysis had referred generally to mathematical methods that employed equations, variables, and constants. With the extensive development of the calculus by Leibniz and his school, analysis became identified with all calculus-related subjects. In addition to this historical association, there was a deeper sense in which analytic methods were fundamental to the new mathematics. An analytic equation implied the existence of a relation that remained valid as the variables changed continuously in magnitude. Analytic algorithms and transformations presupposed a correspondence between local and global change, the basic concern of the calculus. It is this aspect of analysis that fascinated Euler and Lagrange and caused them to see in it the "true metaphysics" of the calculus.

OTHER DEVELOPMENTS

During the period 1600–1800 significant advances occurred in the theory of equations, foundations of Eu-

clidean geometry, number theory, projective geometry, and probability theory. These subjects, which became mature branches of mathematics only in the 19th century, never rivaled analysis and mechanics as programs of research.

Theory of equations. After the dramatic successes of Niccolò Fontana Tartaglia and Lodovico Ferrari in the 16th century the theory of equations developed slowly, as problems resisted solution by known techniques. In the later 18th century the subject experienced an infusion of new ideas. Interest was concentrated on two problems. The first was to establish the existence of a root of the general polynomial equation of degree n . The second was to express the roots as algebraic functions of the coefficients, or to show why it was not in general possible to do so.

The proposition that the general polynomial with real coefficients has a root of the form $a + b\sqrt{-1}$ became known later as the fundamental theorem of algebra. By 1742 Euler recognized that roots appear in conjugate pairs; if $a + b\sqrt{-1}$ is a root, then so is $a - b\sqrt{-1}$. Thus, if $a + b\sqrt{-1}$ is a root of $f(x) = 0$, then $f(x) = (x^2 - 2ax - a^2 - b^2)g(x)$. The fundamental theorem was therefore equivalent to asserting that a polynomial may be decomposed into linear and quadratic factors. This result was of considerable importance for the theory of integration, since by the method of partial fractions it ensured that a rational function, the quotient of two polynomials, could always be integrated in terms of algebraic and elementary transcendental functions.

Although d'Alembert, Euler, and Lagrange worked on the fundamental theorem, the first successful proof was developed by Carl Friedrich Gauss in his doctoral dissertation of 1799. Earlier researchers had investigated special cases or had concentrated on showing that all possible roots were of the form $a \pm b\sqrt{-1}$. Gauss tackled the problem of existence directly. Expressing the unknown in terms of the polar variables r and θ , he showed that a solution of the polynomial would lie at the intersection of two curves of the form $T(r, \theta) = 0$ and $U(r, \theta) = 0$. By a careful and rigorous investigation he proved that the two curves intersect.

Gauss's demonstration of the fundamental theorem initiated a new approach to the question of mathematical existence. In the 18th century mathematicians were interested in the nature of particular analytic processes or the form that given solutions should take. Mathematical entities were regarded as things that were given, not as things whose existence needed to be established. Because analysis was applied in geometry and mechanics, the formalism seemed to possess a clear interpretation that obviated any need to consider questions of existence. Gauss's demonstration was the beginning of a change of attitude in mathematics, of a shift to the rigorous, internal development of the subject.

The problem of expressing the roots of a polynomial as functions of the coefficients was addressed by several mathematicians independently around 1770. The Cambridge mathematician Edward Waring published treatises in 1762 and 1770 on the theory of equations. In 1770 Lagrange presented a long expository memoir on the subject to the Berlin Academy, and in 1771 Alexandre Vandermonde submitted a paper to the French Academy of Sciences. Although the ideas of the three men were related, Lagrange's memoir was the most extensive and most influential historically.

Lagrange presented a detailed analysis of the solution by radicals of second-, third-, and fourth-degree equations and investigated why these solutions failed when the degree was greater than or equal to five. He introduced the novel idea of considering functions of the roots and examining the values they assumed as the roots were permuted. He was able to show that the solution of an equation depends on the construction of a second resolvent equation, but he was unable to provide a general procedure for solving the resolvent when the degree of the original equation was greater than four. Although his theory left the subject in an unfinished condition, it provided a solid basis for future work. The search for a general solution to the polynomial equation would provide the greatest single impetus for the transformation of algebra in the 19th century.

The fundamental theorem of algebra

The theory of functions of a real variable

The efforts of Lagrange, Vandermonde, and Waring illustrate how difficult it was to develop new concepts in algebra. The history of the theory of equations belies the view that mathematics is subject to an almost automatic technical development. Much of the later algebraic work would be devoted to devising terminology, concepts, and methods necessary to advance the subject.

Foundations of geometry. Although the emphasis of mathematics after 1650 was increasingly on analysis, foundational questions in classical geometry continued to arouse interest. Attention centred on the fifth postulate of Book I of the *Elements*, which Euclid had used to prove the existence of a unique parallel through a point to a given line. Since antiquity, Greek, Islāmic, and European geometers had attempted unsuccessfully to show that the parallel postulate need not be a postulate but could instead be deduced from the other postulates of Euclidean geometry. During the period 1600–1800 mathematicians continued these efforts by trying to show that the postulate was equivalent to some result that was considered self-evident. Although the decisive breakthrough to non-Euclidean geometry would not occur until the 19th century, researchers did achieve a deeper and more systematic understanding of the classical properties of space.

Interest in the parallel postulate developed in the 16th century after the recovery and Latin translation of Proclus' commentary on Euclid's *Elements*. The Italian researchers Christopher Clavius in 1574 and Giordano Vitale in 1680 showed that the postulate is equivalent to asserting that the line equidistant from a straight line is a straight line. In 1693 John Wallis, Savilian professor of geometry at Oxford, attempted a different demonstration, proving that the axiom follows from the assumption that to every figure there exists a similar figure of arbitrary magnitude.

In 1733 the Italian Girolamo Saccheri published his *Euclides ab Omni Naevo Vindicatus* (*Euclid Cleared of Every Flaw*). This was an important work of synthesis in which he provided a complete analysis of the problem of parallels in terms of Omar Khayyam's quadrilateral. Using the Euclidean assumption that straight lines do not enclose an area, he was able to exclude geometries that contain no parallels. It remained to prove the existence of a unique parallel through a point to a given line. To do this Saccheri adopted the procedure of *reductio ad absurdum*; he assumed the existence of more than one parallel and attempted to derive a contradiction. After a long and detailed investigation, he was able to convince himself (mistakenly) that he had found the desired contradiction.

In 1766 Johann Heinrich Lambert of the Berlin Academy composed *Die Theorie der Parallelinien* ("The Theory of Parallel Lines"; published 1786), a penetrating study of the fifth postulate in Euclidean geometry. Among other theorems Lambert proved is that the parallel axiom is equivalent to the assertion that the sum of the angles of a triangle is equal to two right angles. He combined this fact with Wallis' result to arrive at an unexpected characterization of classical space. According to Lambert, if the parallel postulate is rejected, it follows that for every angle θ less than $2R/3$ (R is a right angle) an equilateral triangle can be constructed with corner angle θ . By Wallis' result any triangle similar to this triangle must be congruent to it. It is therefore possible to associate with every angle a definite length, the side of the corresponding equilateral triangle. Since the measurement of angles is absolute, independent of any convention concerning the selection of units, it follows that an absolute unit of length exists. Hence, to accept the parallel postulate is to deny the possibility of an absolute concept of length.

The final 18th-century contribution to the theory of parallels was Adrien-Marie Legendre's textbook *Éléments de géométrie*, the first edition of which appeared in 1794. Legendre presented an elegant demonstration that purported to show that the sum of the angles of a triangle is equal to two right angles. He believed that he had conclusively established the validity of the parallel postulate. His work attracted a large audience and was influential in informing readers of the new ideas in geometry.

The 18th-century failure to develop a non-Euclidean geometry was rooted in deeply held philosophical beliefs. In

his *Critique of Pure Reason* (1781), Immanuel Kant had emphasized the synthetic a priori character of mathematical judgments. From this standpoint, statements of geometry and arithmetic were necessarily true propositions with definite empirical content. The existence of similar figures of different size, or the conventional character of units of length, appeared self-evident to mathematicians of the period. As late as 1824 Simon Laplace wrote:

Thus the notion of space includes a special property, self-evident, without which the properties of parallels cannot be rigorously established. The idea of a bounded region, e.g., the circle, contains nothing which depends on its absolute magnitude. But if we imagine its radius to diminish, we are brought without fail to the diminution in the same ratio of its circumference and the sides of all the inscribed figures. This proportionality appears to me a more natural postulate than that of Euclid, and it is worthy of note that it is discovered afresh in the results of the theory of universal gravitation.

(C.G.F.)

Mathematics in the 19th and 20th centuries

Most of the powerful abstract mathematical theories in use today originated in the 19th century, so that any historical account of the period should be supplemented by reference to detailed treatments of these topics. Moreover, mathematics grew so much during this period that any account must necessarily be selective. Nonetheless, some broad features stand out. The growth of mathematics as a profession was accompanied by a sharpening division between mathematics and the physical sciences, and contact between the two subjects takes place today across a clear professional boundary. One result of this separation has been that mathematics, no longer able to rely on its scientific import for its validity, developed markedly higher standards of rigour. It was also freed to develop in directions that had little to do with applicability. Some of these pure creations have turned out to be surprisingly applicable, while the attention to rigour has led to a wholly novel conception of the nature of mathematics and logic. Moreover, many outstanding questions in mathematics yielded to the more conceptual approaches that came into vogue.

Projective geometry. The French Revolution provoked a radical rethinking of education in France, and mathematics was given a prominent role. The École Polytechnique was established in 1794 with the ambitious task of preparing all candidates for the specialist civil and military engineering schools of the republic. Mathematicians of the highest calibre were involved; the result was a rapid and sustained development of the subject. The inspiration for the École was that of Gaspard Monge, who believed strongly that mathematics should serve the scientific and technical needs of the state. To that end he devised a syllabus that promoted his own descriptive geometry, which was useful in the design of forts, gun emplacements, and machines and which was employed to great effect in the Napoleonic survey of Egyptian historical sites.

In Monge's descriptive geometry, three-dimensional objects are described by their orthogonal projections onto a horizontal and a vertical plane, the plan and elevation of the object. A pupil of Monge, Jean-Victor Poncelet, was taken prisoner in Napoleon's retreat from Moscow and sought to keep up his spirits while in jail in Saratov by thinking over the geometry he had learned. He dispensed with the restriction to orthogonal projections and decided to investigate what properties figures have in common with their shadows.

There are several of these properties: a straight line casts a straight shadow, and a tangent to a curve casts a shadow that is tangent to the shadow of the curve. But some properties are lost: the lengths and angles of a figure bear no relation to the lengths and angles of its shadow. Poncelet felt that the properties that survive are worthy of study, and, by considering only those properties that a figure shares with all its shadows, Poncelet hoped to put truly geometric reasoning on a par with algebraic geometry.

In 1822 Poncelet published the *Traité des propriétés projectives des figures* (*Treatise on the Projective Properties of Figures*). From his standpoint, every conic section is

The parallel postulate

The influence of Kant

Poncelet

equivalent to a circle, so his treatise contained a unified treatment of the theory of conic sections. It also established several new results. Geometers who took up his work divided into two groups: those who accepted his terms and those who, finding them obscure, reformulated his ideas in the spirit of algebraic geometry. On the algebraic side it was taken up in Germany by August Ferdinand Möbius, who seems to have come to his ideas independently of Poncelet, and then by Julius Plücker. They showed how rich was the projective geometry of curves defined by algebraic equations and thereby gave an enormous boost to the algebraic study of curves, comparable to the original impetus provided by Descartes. Germany also produced synthetic projective geometers, notably Jakob Steiner (born in Switzerland but educated in Germany) and Karl George Christian von Staudt, who emphasized what can be understood about a figure from a careful consideration of all its transformations.

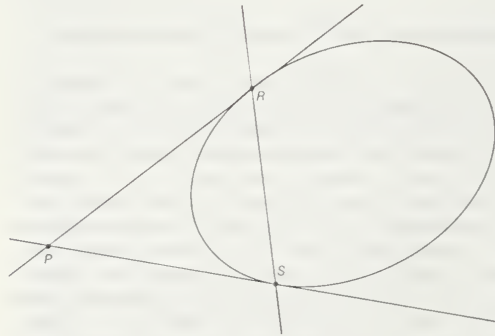


Figure 12: Duality associates with the point P the line RS , and vice versa.

Within the debates about projective geometry emerged one of the few synthetic ideas to be discovered since the days of Euclid, that of duality. This associates with each point a line and with each line a point, in such a way that (1) three points lying in a line give rise to three lines meeting in a point, and, conversely, three lines meeting in a point give rise to three points lying on a line, and (2) if one starts with a point (or a line) and passes to the associated line (point) and then repeats the process, one returns to the original point (line). One way of doing this (presented by Poncelet) is to pick an arbitrary conic and then to associate with a point P lying outside the conic the line that joins the points R and S at which the tangents through P to the conic touch the conic (see Figure 12). A second method is needed for points on or inside the conic. The feature of duality that makes it so exciting is that one can apply it mechanically to every proof in geometry, interchanging "point" and "line" and "collinear" and "concurrent" throughout, and so obtain a new result. Sometimes a result turns out to be equivalent to the original, sometimes to its converse, but at a single stroke the number of theorems was more or less doubled.

Making the calculus rigorous. Monge's educational ideas were opposed by Lagrange, who favoured a more traditional and theoretical diet of advanced calculus and rational mechanics (the application of the calculus to the study of the motion of solids and liquids). Eventually Lagrange won, and the vision of mathematics that was presented to the world was that of an autonomous subject that was also applicable to a broad range of phenomena by virtue of its great generality, a view that has persisted to the present day.

During the 1820s Augustin-Louis Cauchy lectured at the École Polytechnique on the foundations of the calculus. Since its invention it had been generally agreed that the calculus gave correct answers, but no one had been able to give a satisfactory explanation of why this was so. Cauchy rejected Lagrange's algebraic approach and proved that Lagrange's basic assumption that every function had a power series expansion was in fact false. Newton had suggested a geometric or dynamic basis for calculus, but this ran the risk of introducing a vicious circle when the calculus was applied to mechanical or geometric problems. Cauchy

proposed basing the calculus on a sophisticated and difficult interpretation of the idea of two points or numbers being arbitrarily close together. Although his students disliked the new approach, and Cauchy was ordered to teach material that the students could actually understand and use, his methods gradually became established and refined to form the core of the modern rigorous calculus, a subject now called mathematical analysis.

Traditionally, the calculus had been concerned with the two processes of differentiation and integration and the reciprocal relation that exists between them. Cauchy provided a novel underpinning by stressing the importance of the concept of continuity, which is more basic than either. He showed that, once the concepts of a continuous function and limit are defined, the concepts of a differentiable function and an integrable function can be defined in terms of them. Unfortunately, neither of these concepts is easy to grasp, and the much-needed degree of precision they bring to mathematics has proved difficult to appreciate. Roughly speaking, a function is continuous at a point in its domain if small changes in the input around the specified value only produce small changes in the output.

Thus, in Figure 13 the familiar graph of a parabola $y = x^2$ is continuous around the point $x = 0$; as x varies by small amounts, so necessarily does y . On the other hand, the graph of the function that takes the value 0 when x is negative or zero, and the value 1 when x is positive, plainly has a discontinuous graph at the point $x = 0$, and it is indeed discontinuous there according to the definition. If x varies from 0 by any small positive amount, the value of the function jumps by the fixed amount 1, which is not an arbitrarily small amount.

Cauchy said that a function $f(x)$ tends to a limiting value 1 as x tends to the value a whenever the value of the difference $f(x) - f(a)$ becomes arbitrarily small as the difference $x - a$ itself becomes arbitrarily small. He then showed that, if $f(x)$ is continuous at a , the limiting value of the function as x tended to a was indeed $f(a)$. The crucial feature of this definition is that it defines what it means for a variable quantity to tend to something entirely without reference to ideas of motion.

Cauchy then said a function $f(x)$ is differentiable at the point a if, as x tends to a (which it is never allowed to reach), the value of the quotient $\{f(x) - f(a)\}/(x - a)$ tends to a limiting value, called the derivative of the function $f(x)$ at a . To define the integral of a function $f(x)$ between the values a and b , Cauchy went back to the primitive idea of integral as the measure of the area under the graph of the function. He approximated this area by rectangles and said that, if the sum of the areas of the rectangles tends to a limit as their number increases indefinitely and

Mathematical analysis

Definition of the integral

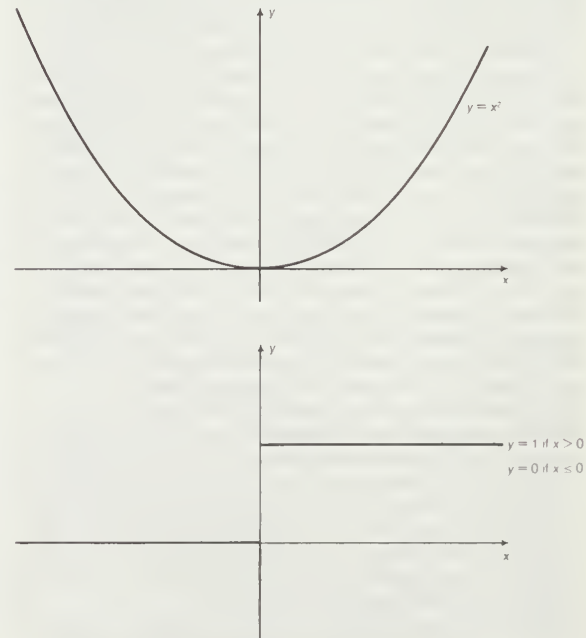


Figure 13: Continuous and discontinuous functions.

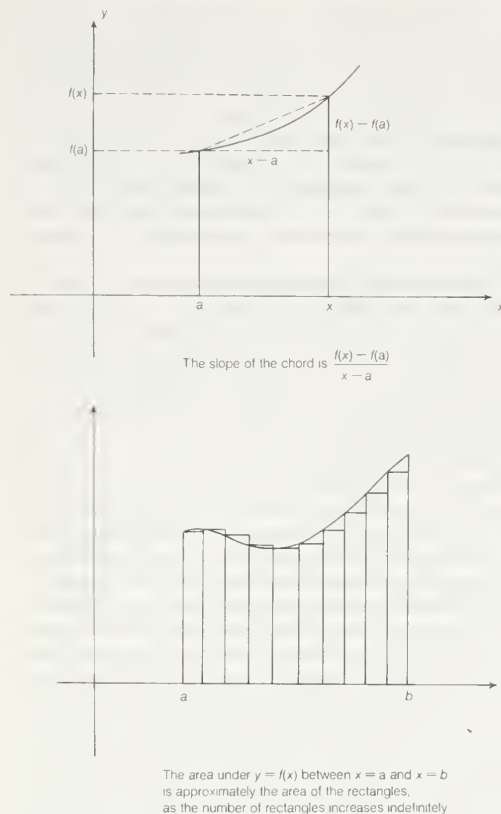


Figure 14: Differentiation and integration.

if this limiting value is the same however the rectangles are obtained, then the function is integrable. Its integral is the common limiting value. After he had defined the integral independently of the differential calculus, Cauchy had to prove that the processes of integrating and differentiating are mutually inverse. This he did, giving for the first time a rigorous foundation to all the elementary calculus of his day.

Fourier series. The other crucial figure of the time in France was Joseph Fourier. His major contribution, presented in *La Théorie analytique de la chaleur* (1822: *The Analytical Theory of Heat*), was to the theory of heat diffusion in solid bodies. He proposed that any function could be written as an infinite sum of the trigonometric functions cosine and sine, for example,

$$f(x) = a_0 + a_1 \sin x + a_2 \sin 2x + \dots$$

While expressions of this kind had been written down earlier, what was new in Fourier's treatment was the degree of attention given to their convergence. He investigated the question, "Given the function $f(x)$, for what range of values of x does the expression above sum to a finite number?" It turned out that the answer depends on the coefficients a_n , and Fourier gave rules for obtaining them of the form:

$$a_n = \int_{-\pi}^{\pi} f(x) \sin (nx) dx.$$

Had Fourier's work been entirely correct, it would have brought all functions into the calculus, making possible the solution of many kinds of differential equations and greatly extending the theory of mathematical physics. But his arguments were unduly naive: after Cauchy, it was not clear that the function $f(x) \sin (nx)$ was necessarily integrable. When Fourier's ideas were finally published, they were eagerly taken up, but the more cautious mathematicians, notably the influential German mathematician Peter Gustav Lejeune Dirichlet, wanted to rederive Fourier's conclusions in a more rigorous way. Fourier's methodology was widely accepted, but questions about its validity in detail were to occupy mathematicians for the rest of the century.

Elliptic functions. The theory of functions of a complex variable was also being decisively reformulated. At

the start of the 19th century, complex numbers were discussed from a quasi-philosophical standpoint by several French writers, notably Jean Robert Argand. A consensus emerged that complex numbers should be thought of as pairs of real numbers, with suitable rules for their addition and multiplication so that the pair $(0, 1)$ was a square root of -1 . The underlying meaning of such a number pair was given by its geometric interpretation either as a point in a plane or as a directed segment joining the coordinate origin to the point in question. (This representation is sometimes called the Argand diagram.) In 1827, while revising an earlier manuscript for publication, Cauchy showed how the problem of integrating functions of two variables can be illuminated by a theory of functions of a single complex variable, which he was then developing. But the decisive influence on the growth of the subject came from the theory of elliptic functions.

The study of elliptic functions originated in the 18th century, when many authors studied integrals of the form

$$\int_0^x \frac{p(t) dt}{\sqrt{q(t)}}$$

where $p(t)$ and $q(t)$ are polynomials in t and $q(t)$ is of degree 3 or 4 in t . Such integrals arise naturally, for example, as an expression for the length of an arc of an ellipse (whence the name). These integrals cannot be evaluated explicitly; they do not define a function that can be obtained from the rational and trigonometric functions, a difficulty that added to their interest. Elliptic integrals were intensively studied for many years by the French mathematician Legendre, who was able to calculate tables of values for such expressions as functions of their upper endpoint, x . But the topic was completely transformed in the late 1820s by the independent but closely overlapping discoveries of two young mathematicians, the Norwegian Niels Henrik Abel and the German Carl Gustav Jacob Jacobi. These men showed that, if one allowed the variable x to be complex and the problem were inverted, so that the object of study became

$$u = \int_0^x \frac{p(t) dt}{\sqrt{q(t)}}$$

considered as defining a function x of a variable u , then a remarkable new theory could be discovered. The new function, for example, possessed a property that generalized the basic property of periodicity of the trigonometric functions sine and cosine: $\sin (x) = \sin (x + 2\pi)$. Any function of the kind just described has two distinct periods, ω_1 and ω_2 :

$$x(u) = x(u + \omega_1) = x(u + \omega_2).$$

These new functions, the elliptic functions, aroused a considerable degree of interest from the first because many interesting and novel things could be said about them. The analogy with trigonometric functions ran very deep (indeed the trigonometric functions turned out to be special cases of elliptic functions), but their greatest influence was on the burgeoning general study of functions of a complex variable. The theory of elliptic functions became the paradigm of what could be discovered by allowing variables to be complex instead of real. But their natural generalization to functions defined by more complicated

The Argand diagram

Elliptic functions



Figure 15: The plane of complex numbers (Argand diagram).

integrands, although it yielded partial results, resisted analysis until the second half of the 19th century.

The theory of numbers. While the theory of elliptic functions typifies the 19th century's enthusiasm for pure mathematics, some contemporary mathematicians said that the simultaneous developments in number theory carried that enthusiasm to excess. Nonetheless, during the 19th century the algebraic theory of numbers grew from being a minority interest to its present central importance in pure mathematics. The earlier investigations of Fermat had eventually drawn the attention of Euler and Lagrange. Euler proved some of Fermat's unproved claims and discovered many new and surprising facts; Lagrange not only supplied proofs of many remarks that Euler had merely conjectured but also worked them into something like a coherent theory. For example, it was known to Fermat that the numbers which can be written as the sum of two squares are either the number 2, squares themselves, primes of the form $4n + 1$, or else products of these numbers. Thus 29, which is $4 \times 7 + 1$, is $5^2 + 2^2$, but 35, which is not of this form, cannot be written as the sum of two squares. Euler had proved this result and had gone on to consider similar cases, such as primes of the form $x^2 + 2y^2$, or $x^2 + 3y^2$. But it was left to Lagrange to provide a general theory covering all expressions of the form $ax^2 + bxy + cy^2$, quadratic forms, as they are called.

Lagrange's theory of quadratic forms had made considerable use of the idea that a given quadratic form could often be simplified to another with the same properties but with smaller coefficients. To do this in practice, it was often necessary to consider whether a given integer left a remainder that was a square when it was divided by another given integer. (For example, 48 leaves a remainder of 4 upon division by 11, and 4 is a square.) Legendre discovered a remarkable connection between the question, "Does the integer p leave a square remainder on division by q ?" and the seemingly unrelated question, "Does the integer q leave a square remainder upon division by p ?" He saw, in fact, that, when p and q are primes, both questions have the same answer unless both primes are of the form $4n - 1$. Because this observation connects two questions in which the integers p and q play mutually opposite roles, it became known as the law of quadratic reciprocity. Legendre also gave an effective way of extending his law to cases when p and q are not prime.

All this work set the scene for the emergence of Gauss, whose *Disquisitiones Arithmeticae* not only consummated what had gone before but also directed number theorists in new and deeper directions. He rightly showed that Legendre's proof of the law of quadratic reciprocity was fundamentally flawed and gave the first rigorous proof. His work suggested that there were profound connections between the original question and other branches of number theory, a fact that he perceived to be of signal importance for the subject. He extended Lagrange's theory of quadratic forms by showing how two quadratic forms can be "multiplied" to obtain a third. Later mathematicians were to rework this into an important example of the theory of finite commutative groups. And in the long final section of his book Gauss gave the theory that lay behind his first discovery as a mathematician: that a regular 17-sided figure can be constructed by circle and straightedge alone.

The discovery that the regular "17-gon" is so constructible was the first such discovery since the Greeks, who had known only of the equilateral triangle, the square, the regular pentagon, the regular 15-sided figure, and the figures that can be obtained from these by successively bisecting all the sides. But what was of much greater significance than the discovery was the theory that underpinned it, the theory of what are now called algebraic numbers. It may be thought of as an analysis of how complicated a number may be while yet being amenable to an exact treatment.

The simplest numbers to understand and use are the integers and the rational numbers; the irrational numbers seem to pose problems. Famous among these is $\sqrt{2}$. It cannot be written as a finite or repeating decimal (because it is not rational), but it can be manipulated algebraically very easily. It is only necessary to replace every occur-

rence of $(\sqrt{2})^2$ by 2. In this way expressions of the form $m + n\sqrt{2}$, where m and n are integers, can be handled arithmetically. These expressions have many properties akin to those of whole numbers, and one can even define prime numbers of this form; therefore they are called algebraic integers. In this case they are obtained by grafting onto the rational numbers a solution of the polynomial equation $x^2 - 2 = 0$. In general an algebraic integer is any solution, real or complex, of a polynomial equation with integer coefficients in which the coefficient of the highest power of the unknown is 1.

Gauss's theory of algebraic integers led to the question of determining when a polynomial of degree n with integer coefficients can be solved given the solvability of polynomial equations of lower degree but with coefficients that are algebraic integers. For example, he regarded the coordinates of the 17 vertices of a regular 17-sided figure as complex numbers satisfying the equation $x^{17} - 1 = 0$, and thus as algebraic integers. One such integer is $z = 1$. He showed that the rest are obtained by solving a succession of four quadratic equations. Because solving a quadratic equation is equivalent to performing a construction with a ruler and compass, as Descartes had shown long before, Gauss had shown how to construct the regular 17-gon.

Inspired by Gauss's works on the theory of numbers, a growing school of mathematicians was drawn to the subject. Like Gauss, the German mathematician Ernst Eduard Kummer sought to generalize the law of quadratic reciprocity to deal with questions about third, fourth, and higher powers of numbers. He found that his work led him in an unexpected direction, toward a partial resolution of Fermat's last theorem. In 1637 Fermat wrote in the margin of his copy of Diophantus' *Arithmetica* the claim to have a proof that there are no solutions in positive integers to the equation $x^n + y^n = z^n$ if $n > 2$. However, no proof was ever discovered among his notebooks.

Kummer's approach was to develop the theory of algebraic integers. If it could be shown that the equation had no solution in suitable algebraic integers, then a fortiori there could be no solution in ordinary integers. He was eventually able to establish the truth of Fermat's last theorem for a large class of prime exponents n (those satisfying some technical conditions needed to make the proof work). This was the first significant breakthrough in the study of the theorem. Together with the earlier work of the French mathematician Sophie Germain, it has enabled mathematicians to establish Fermat's last theorem for every value of n from 3 to 4,000,000. However, Kummer's way around the difficulties he encountered further propelled the theory of algebraic integers into the realm of abstraction. It amounted to the suggestion that there should be yet other types of integers, but many found these ideas obscure.

In Germany Richard Dedekind patiently created a new approach, in which each new number (called ideal) was defined by means of a suitable set of algebraic integers in such a way that it was the common divisor of the set of algebraic integers used to define it. Dedekind's work was slow to gain approval, yet it illustrates several of the most profound features of modern mathematics. It was clear to Dedekind that the ideal algebraic integers were the work of the human mind. Their existence can neither be based on nor deduced from the existence of physical objects, analogies with natural processes, or some process of abstraction from more familiar things. A second feature of Dedekind's work was its reliance on the idea of sets of objects, such as sets of numbers, even sets of sets. Dedekind's work showed how basic the naive conception of a set could be. The third crucial feature of his work was its emphasis on the structural aspects of algebra. The presentation of number theory as a theory about objects that can be manipulated (in this case, added and multiplied) according to certain rules akin to those governing ordinary numbers was to be a paradigm of the more formal theories of the 20th century.

The theory of equations. Another subject that was transformed in the 19th century was the theory of equations. Ever since Tartaglia and Ferrari in the 16th century had found rules giving the solutions of cubic and quartic

The theory of algebraic numbers

The ideal algebraic integers

equations in terms of the coefficients of the equations, formulas had unsuccessfully been sought for equations of the fifth and higher degrees. At stake was the existence of a formula that expresses the roots of a quintic equation in terms of the coefficients. This formula moreover, must involve only the operations of addition, subtraction, multiplication, and division, together with the extraction of roots, since that was all that had been required for the solution of quadratic, cubic, and quartic equations. If such a formula were to exist, the quintic would accordingly be said to be solvable by radicals.

In 1770 Lagrange had analyzed all the successful methods he knew for equations of degrees 2, 3, and 4, in an attempt to see why they worked and how they could be generalized. His analysis of the problem in terms of permutations of the roots was promising, but he became more and more doubtful as the years went by that his complicated line of attack could be carried through. The first valid proof that the general quintic is not solvable by radicals was offered only after his death, in a startlingly short paper by Abel, written in 1824.

Abel also showed by example that some quintic equations were solvable by radicals and that some equations could be solved unexpectedly easily. For example, the equation $x^5 - 1 = 0$ has one root $x = 1$, but the remaining four roots can be found just by extracting square roots, not fourth roots as might be expected. He therefore raised the question, "What equations of degree higher than 4 are solvable by radicals?"

Abel died in 1829 at the age of 26 and did not resolve the problem he had posed. Almost at once, however, the astonishing prodigy Évariste Galois burst upon the Parisian mathematical scene. He submitted an account of his novel theory of equations to the Academy of Sciences in 1829, but the manuscript was lost. A second version was also lost and was not found among Fourier's papers when Fourier, the secretary of the academy, died in 1830. Galois was killed in a duel in 1832, at the age of 20, and it was not until his papers were published in Joseph Liouville's *Journal de mathématiques* in 1846 that his work began to receive the attention it deserved. His theory eventually made the theory of equations into a mere part of the theory of groups. Galois emphasized the group (as he called it) of permutations of the roots of an equation. This move took him away from the equations themselves, instead turning toward the markedly more tractable study of permutations. To any given equation there corresponds a definite group, with a definite collection of subgroups. To explain which equations were solvable by radicals and which were not, Galois analyzed the ways in which these subgroups were related to one another: solvable equations gave rise to what are now called a chain of normal subgroups with cyclic quotients. This technical condition makes it clear how far mathematicians had gone from the familiar questions of 18th-century mathematics, and it marks a transition characteristic of modern mathematics: the replacement of formal calculation by conceptual analysis. This is a luxury available to the pure mathematician that the applied mathematician faced with a concrete problem cannot always afford.

According to this theory, a group is a set of objects that one can combine in pairs in such a way that the resulting object is also in the set. Moreover, this way of combination has to obey the following rules (here objects in the group are denoted a, b , etc., and the combination of a and b is written $a * b$):

1. There is an element, e , such that $a * e = a = e * a$ for every element a in the group. This element is called the identity element of the group.
2. For every element a there is an element, written a^{-1} , with the property that $a * a^{-1} = e = a^{-1} * a$. The element a^{-1} is called the inverse of a .
3. For every a, b , and c in the group the associative law holds: $(a * b) * c = a * (b * c)$.

Examples of groups include the integers with $*$ interpreted as addition and the positive rational numbers with $*$ interpreted as multiplication. An important property shared by some groups but not all is commutativity: for every element a and b , $a * b = b * a$. The rotations of an

object in the plane around a fixed point form a commutative group, but the rotations of a three-dimensional object around a fixed point form a noncommutative group.

Gauss. A convenient way to assess the situation in mathematics in the mid-19th century is to look at the career of its greatest exponent, the last man to be called the "Prince of Mathematics," Carl Friedrich Gauss. In 1801, the same year in which he published his *Disquisitiones Arithmeticae*, he rediscovered the asteroid Ceres (which had disappeared behind the Sun not long after it had first been discovered and before its orbit was precisely known). He was the first to give a sound analysis of the method of least squares in the analysis of statistical data. Gauss did important work in potential theory and, with the German physicist Wilhelm Weber, built the first electric telegraph. He helped conduct the first survey of the Earth's magnetic field and did both theoretical and field work in cartography and surveying. He was a polymath who almost single-handedly embraced what elsewhere was being put asunder: the world of science and the world of mathematics. It is his purely mathematical work, however, that in its day and ever since has been regarded as the best evidence of his genius.

Gauss's writings transformed the theory of numbers. His theory of algebraic integers lay close to the theory of equations as Galois was to redefine it. More remarkable are the extensive writings on the theory of elliptic functions dating from 1797 to the 1820s but unpublished at his death. In 1827 he published his crucial discovery that the curvature of a surface can be defined intrinsically—that is, solely in terms of properties defined within the surface and without reference to the surrounding Euclidean space. This result was to be decisive in the acceptance of non-Euclidean geometry. All his work displays a sharp concern for rigour and a refusal to rely on intuition or physical analogy, which was to serve as an inspiration to his successors. His emphasis on achieving full conceptual understanding, which may have led to his dislike of publication, was by no means the least influential of his achievements.

Gauss's work on the theory of elliptic functions

The theory of groups

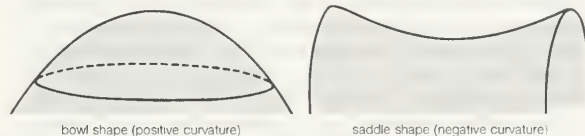


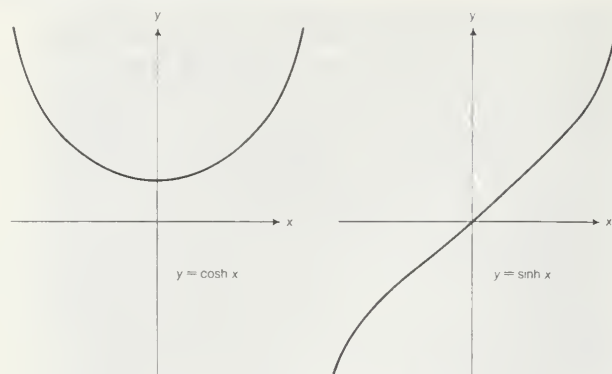
Figure 16: Intrinsic curvature of a surface.

Non-Euclidean geometry. Perhaps it was this desire for conceptual understanding that made Gauss reluctant to publish the fact that he was led more and more "to doubt the truth of geometry," as he put it. For if there was a logically consistent geometry differing from Euclid's only because it made a different assumption about the behaviour of parallel lines, it too could apply to physical space, and so the truth of (Euclidean) geometry could no longer be assured a priori, as Kant had thought.

Gauss's investigations into the new geometry went further than any one else's before him, but he did not publish them. The honour of being the first to proclaim the existence of a new geometry belongs to two others, who did so in the late 1820s: Nicolay Ivanovich Lobachevsky in Russia and János Bolyai in Hungary. Because the similarities in the work of these two men far exceed the differences, it is convenient to describe their work together.

Both men made an assumption about parallel lines that differed from Euclid's and proceeded to draw out its consequences. This way of working cannot guarantee the consistency of one's findings, so strictly speaking they could not prove the existence of a new geometry in this way. Both men described a three-dimensional space different from Euclidean space, couching their findings in the language of trigonometry. The formulas they obtained were exact analogues of the formulas that describe triangles drawn on the surface of a sphere, with the usual trigonometric functions replaced by those of hyperbolic trigonometry. (The functions hyperbolic cosine, written cosh, and hyperbolic sine, written sinh [see Figure 17], are defined as follows: $\cosh x = \frac{1}{2}(e^x + e^{-x})$, and $\sinh x = \frac{1}{2}(e^x - e^{-x})$. They are called hyperbolic because of their use in describ-

The hyperbolic functions

Figure 17: The hyperbolic functions $\cosh x$ and $\sinh x$.

ing the hyperbola. Their names derive from the evident analogy with the trigonometric functions, which Euler showed satisfy these equations: $\cos x = \frac{1}{2}(e^{ix} + e^{-ix})$, and $\sin x = \frac{1}{2i}(e^{ix} - e^{-ix})$. The formulas were what gave their work the precision needed to give conviction in the absence of a sound logical structure. Both men observed that it had become an empirical matter to determine the nature of space, Lobachevsky even going so far as to conduct astronomical observations, although these proved inconclusive.

The work of Bolyai and Lobachevsky was poorly received. Gauss endorsed what they had done, but so discreetly that most mathematicians did not find out his true opinion on the subject until he was dead. The main obstacle each man faced was surely the shocking nature of their discovery. It was easier, and in keeping with 2,000 years of tradition, to continue to believe that Euclidean geometry was correct and that Bolyai and Lobachevsky had somewhere gone astray, like many an investigator before them.

The turn toward acceptance came in the 1860s after Bolyai and Lobachevsky had died. The Italian mathematician Eugenio Beltrami decided to investigate Lobachevsky's work and to place it, if possible, within the context of differential geometry as redefined by Gauss. He therefore moved independently in the direction already taken by Bernhard Riemann. He investigated the surface of constant negative curvature and found that on such a surface triangles obeyed the formulas of hyperbolic trigonometry that Lobachevsky had discovered were appropriate to his form of non-Euclidean geometry. Thus he gave the first rigorous description of a geometry other than Euclid's. Beltrami's account of the surface of constant negative curvature was ingenious. He said that it was an abstract surface which he could describe by drawing maps of it, much as one might describe a sphere by means of the pages of a geographic atlas. He did not claim to have constructed the surface embedded in Euclidean three-dimensional space; later Hilbert showed that it cannot be done.

Riemann. When Gauss died in 1855, his post at Göttingen was taken by Peter Gustav Lejeune Dirichlet. One mathematician who found the presence of Dirichlet a stimulus to research was Riemann, and his few short contributions to mathematics were among the most influential of the century. Riemann's first paper, his doctoral thesis (1851) on the theory of complex functions, provided the foundations for a geometric treatment of functions of a complex variable. His main result guaranteed the existence of a wide class of complex functions satisfying only modest general requirements and so made it clear that complex functions could be expected to occur widely in mathematics. More importantly, it was achieved by yoking together the complex theory with the theory of harmonic functions and potential theory. The theories of complex and harmonic functions were henceforth inseparable.

Riemann then wrote on the theory of Fourier series and their integrability. His paper was directly in the tradition that ran from Cauchy and Fourier to Dirichlet, and it marked a considerable step forward in the precision with which the concept of integral can be defined. In 1854 he took up a subject that much interested Gauss, the hypotheses lying at the basis of geometry.

The study of geometry has always been one of the central concerns of mathematicians. It was the language, and the principal subject matter, of Greek mathematics, the mainstay of elementary education in the subject, and it has an obvious visual appeal. It seems easy to apply, for one can proceed from a base of naively intelligible concepts. In keeping with the general trends of the century, however, it was just the naive concepts that Riemann chose to refine. What he proposed as the basis of geometry was far more radical and fundamental than anything that had gone before.

Riemann took his inspiration from Gauss's discovery that the curvature of a surface is intrinsic, and he argued that one should therefore ignore Euclidean space and treat each surface by itself. A geometric property, he argued, was one that was intrinsic to the surface. To do geometry, it was enough to be given a set of points and a way of measuring lengths along curves in the surface. For this, traditional ways of applying the calculus to the study of curves could be made to suffice. But Riemann did not stop with surfaces. He proposed that geometers study spaces of any dimension in this spirit, even, he said, spaces of infinite dimension.

Several profound consequences followed from this view. It dethroned Euclidean geometry, which now became just one of many geometries. It allowed the geometry of Bolyai and Lobachevsky to be recognized as the geometry of a surface of constant negative curvature, thus resolving doubts about the logical consistency of their work. It highlighted the importance of intrinsic concepts in geometry. It helped open the way to the study of spaces of many dimensions. Last, but not least, Riemann's work ensured that any investigation of the geometric nature of physical space would thereafter have to be partly empirical. One could no longer say that physical space is Euclidean because there is no geometry but Euclid's. This finally destroyed any hope that questions about the world could be answered by a priori reasoning.

In 1857 Riemann published several papers applying his very general methods for the study of complex functions to various parts of mathematics. One of these papers solved the outstanding problem of extending the theory of elliptic functions to the integration of any algebraic function. It opened up the theory of complex functions of several variables and showed how Riemann's novel topological ideas were essential in the study of complex functions. (In subsequent lectures Riemann showed how the special case of the theory of elliptic functions could be regarded as the study of complex functions on a torus.)

Another paper dealt with the question of how many prime numbers there are that are less than any given number x . The answer is a function of x , and Gauss had conjectured on the basis of extensive numerical evidence that this function was approximately $x/\ln(x)$. This turned out to be true, but it was not proved until 1896, when both the Belgian mathematician Charles Jean de la Vallée-Poussin and the French mathematician Jacques-Salomon Hadamard independently proved it. It is remarkable that a question about integers led to a discussion of functions of a complex variable, but similar connections had previously been made by Dirichlet. Riemann took the expression $\Pi(1 - p^{-s})^{-1} = \sum n^{-s}$, introduced by Euler the century before, where the infinite product is taken over all prime numbers p and the sum over all whole numbers n and treated it as a function of s . The infinite sum makes sense whenever s is real and greater than 1. Riemann proceeded to study this function when s is complex (now called the Riemann zeta function), and he thereby not only helped clarify the question of the distribution of primes but also was led to several other remarks that later mathematicians were to find of exceptional interest. One remark has continued to elude proof and remains one of the greatest conjectures in mathematics: the claim that the nonreal zeros of the zeta function are complex numbers whose real part is always equal to $1/2$.

Riemann's influence. In 1859 Dirichlet died and Riemann became a full professor, but he was already ill with tuberculosis, and in 1862 his health broke. He died in 1866. His work, however, exercised a growing influence

Riemann's work on the foundations of geometry

The Riemann zeta function

on his successors. His work on trigonometric series, for example, led to a deepening investigation of the question of when a function is integrable. Attention was concentrated on the nature of the sets of points at which functions and their integrals (when these existed) had unexpected properties. The conclusions that emerged were at first obscure, but it became clear that some properties of point sets were important in the theory of integration, while others were not. (These other properties proved to be a vital part of the emerging subject of topology.) The properties of point sets that matter in integration have to do with the size of the set. If one can change the values of a function on a set of points without changing its integral, it is said that the set is of negligible size. The naive idea is that integrating is a generalization of counting: negligible sets do not need to be counted. Around the turn of the century the French mathematician Henri-Léon Lebesgue managed to systematize this naive idea into a new theory about the size of sets, which included integration as a special case. In this theory, called measure theory, there are sets that can be measured, and they either have positive measure or are negligible (they have zero measure), and there are sets that cannot be measured at all.

The first success for Lebesgue's theory was that, unlike the Cauchy-Riemann integral, it obeyed the rule that, if a sequence of functions $f_n(x)$ tended suitably to a function $f(x)$, then the sequence of integrals $\int f_n(x)dx$ tended to the integral $\int f(x)dx$. This has made it the natural theory of the integral when dealing with questions about trigonometric series. Another advantage is that it is very general. For example, in probability theory it is desirable to estimate the likelihood of certain outcomes of an experiment. By imposing a measure on the space of all possible outcomes, the Russian mathematician Andrey Kolmogorov was the first to put probability theory on a rigorous mathematical footing.

Another example is provided by a remarkable result discovered by the 20th-century American mathematician Norbert Wiener: within the set of all continuous functions on an interval, the set of differentiable functions has measure zero. In probabilistic terms, therefore, the chance that a function taken at random is differentiable has probability zero. In physical terms, this means that, for example, a particle moving under Brownian motion almost certainly is moving on a nondifferentiable path. This discovery clarified Einstein's fundamental ideas about Brownian motion (displayed by the continual motion of specks of dust in a fluid under the constant bombardment of surrounding molecules). The hope of physicists is that Richard Feynman's theory of quantum electrodynamics will yield to a similar measure-theoretic treatment, for it has the disturbing aspect of a theory that has not been made rigorous mathematically but that accords excellently with observation.

Yet another setting for Lebesgue's ideas was to be the theory of Lie groups. The Hungarian mathematician Alfréd Haar showed how to define the concept of measure so that functions defined on Lie groups could be integrated. This became a crucial part of Hermann Weyl's way of representing a Lie group as acting linearly on the space of all (suitable) functions on the group (for technical reasons, "suitable" means that the square of the function is integrable with respect to a Haar measure on the group).

Differential equations. Another field that developed considerably in the 19th century was the theory of differential equations. The pioneer in this direction once again was Cauchy. Above all, he insisted that one should prove that solutions do indeed exist; it is not a priori obvious that every ordinary differential equation has solutions. The methods that Cauchy proposed for these problems fitted naturally into his program of providing rigorous foundations for all the calculus. The solution method he preferred, although the less general of his two approaches, worked equally well in the real and complex cases. It established the existence of a solution equal to the one obtainable by traditional power series methods using newly developed techniques in his theory of functions of a complex variable.

The harder part of the theory of differential equations

concerns partial differential equations, those for which the unknown function is a function of several variables. In the early 19th century there was no known method of proving that a given second- or higher-order partial differential equation had a solution, and there was not even a method of writing down a plausible candidate. In this case progress was to be much less marked. Cauchy found new and more rigorous methods for first-order partial differential equations, but the general case eluded treatment.

An important special case was successfully prosecuted, that of dynamics. Dynamics is the study of the motion of a physical system under the action of forces. Working independently of each other, William Rowan Hamilton in Ireland and Jacobi in Germany showed how problems in dynamics could be reduced to systems of first-order partial differential equations. From this base grew an extensive study of certain partial differential operators. These are straightforward generalizations of a single partial differentiation ($\partial/\partial x$) to a sum of the form

$$a_1 \frac{\partial}{\partial x_1} + \dots + a_n \frac{\partial}{\partial x_n},$$

where the a 's are functions of the x 's. The effect of performing several of these in succession can be complicated, but Jacobi and the other pioneers in this field found that there are formal rules which such operators tend to satisfy. This enabled them to shift attention to these formal rules, and gradually an algebraic analysis of this branch of mathematics began to emerge.

The most influential worker in this direction was the Norwegian, Sophus Lie. Lie, and independently Wilhelm Killing in Germany, came to suspect that the systems of partial differential operators they were studying came in a limited variety of types. Once the number of independent variables was specified (which fixed the dimension of the system), a large class of examples, including many of considerable geometric significance, seemed to fall into a small number of patterns. This suggested that the systems could be classified, and such a prospect naturally excited mathematicians. After much work by Lie and Killing and later by the French mathematician Élie-Joseph Cartan, they were classified. Initially, this discovery aroused interest because it produced order where previously the complexity had threatened chaos and because it could be made to make sense geometrically. The realization that there were to be major implications of this work for the study of physics lay well in the future.

Linear algebra. Differential equations, whether ordinary or partial, may profitably be classified as linear or nonlinear; linear differential equations are those for which the sum of two solutions is again a solution. The equation giving the shape of a vibrating string is linear, which provides the mathematical reason why a string may simultaneously emit more than one frequency. The linearity of an equation makes it easy to find all its solutions, so in general linear problems have been tackled successfully, while nonlinear equations continue to be difficult. Indeed, in many linear problems there can be found a finite family of solutions with the property that any solution is a sum of them (suitably multiplied by arbitrary constants). Obtaining such a family, called a basis, and putting them into their simplest and most useful form, was an important source of many techniques in the field of linear algebra.

Consider, for example, the system of linear differential equations

$$\frac{dy_1}{dx} = ay_1 + by_2, \quad \frac{dy_2}{dx} = cy_1 + dy_2.$$

It is evidently much more difficult to study than the system $dy_1/dx = ay_1$, $dy_2/dx = by_2$, whose solutions are (constant multiples of) $y_1 = \exp(ax)$ and $y_2 = \exp(bx)$. But if a suitable linear combination of y_1 and y_2 can be found so that the first system reduces to the second, then it is enough to solve the second system. The existence of such a reduction is determined by the array (called a matrix) of the four numbers

$\begin{bmatrix} a & b \\ c & d \end{bmatrix}$. In 1858 the English mathematician Arthur

The classification of systems of partial differential operators

Cayley began the study of matrices in their own right when he noticed that they satisfy polynomial equations. The

matrix $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, for example, satisfies the equation

$A^2 - (a+d)A + (ad-bc) = 0$. Moreover, if this equation has two distinct roots, say, α and β , then the sought-for reduction will exist and the coefficients of the simpler system will indeed be those roots α and β . If the equation has a repeated root, then the reduction usually cannot be carried out. In either case the difficult part of solving the original differential equation has been reduced to elementary algebra.

Vector spaces

The study of linear algebra begun by Cayley and continued by Leopold Kronecker includes a powerful theory of vector spaces. These are sets whose elements can be added together and multiplied by arbitrary numbers, such as the family of solutions of a linear differential equation. A more familiar example is that of three-dimensional space. If one picks an origin, then every point in space can be labeled by the line segment (called a "vector") joining it to the origin. Matrices appear as ways of representing linear transformations of a vector space—i.e., transformations that preserve sums and multiplication by numbers: the transformation T is linear if, for any vectors u, v , $T(u+v) = T(u) + T(v)$ and, for any scalar λ , $T(\lambda v) = \lambda T(v)$. When the vector space is finite-dimensional, linear algebra and geometry form a potent combination. Vector spaces of infinite dimensions also are studied.

The theory of vector spaces is useful in other ways. Vectors in three-dimensional space represent such physically important concepts as velocities and forces. Such an assignment of vector to point is called a vector field; examples include electric and magnetic fields. Scientists such as Maxwell and J. Willard Gibbs took up vector analysis and were able to extend vector methods to the calculus. They introduced in this way measures of how a vector field varies infinitesimally, which, under the names div, grad, and curl, have become the standard tools in the study of electromagnetism and potential theory. To the modern mathematician, div, grad, and curl form part of a theory to which Stokes's theorem (a special case of which is Green's theorem) is central. This result, named after two leading English applied mathematicians of the 19th century, generalizes the fundamental theorem of the calculus to functions of several variables. The fundamental theorem asserts that

$$\int_a^b f'(x) = f(b) - f(a),$$

which can be read as saying that the integral of the derivative of some function in an interval is equal to the difference in the values of the function at the endpoints of the interval. Generalized to a part of a surface or space, this asserts that the integral of the derivative of some function over a region is equal to the integral of the function over the boundary of the region. In symbols this says that $\int d\omega = \int \omega$, where the first integral is taken over the region in question and the second integral over its boundary, while $d\omega$ is the derivative of ω .

The foundations of geometry. By the late 19th century the hegemony of Euclidean geometry had been challenged by non-Euclidean geometry and projective geometry. The first notable attempt to reorganize the study of geometry was made by the German mathematician Felix Klein and published at Erlangen in 1872. In his *Erlanger Programm* Klein proposed that Euclidean and non-Euclidean geometry be regarded as special cases of projective geometry. In each case the common features that, in Klein's opinion, made them geometries were that there was a set of points, called a "space," and a group of transformations by means of which figures could be moved around in the space without altering their essential properties. For example, in Euclidean plane geometry the space is the familiar plane, and the transformations are rotations, reflections, translations, and their composites, none of which change either length or angle, the basic properties of figures in Euclidean geometry. Different geometries would have dif-

ferent spaces and different groups, and the figures would have different basic properties.

Klein produced an account that unified a large class of geometries, roughly speaking all those which were homogeneous in the sense that every piece of the space looked like every other piece of the space. This excluded, for example, geometries on surfaces of variable curvature, but it produced an attractive package for the rest and gratified the intuition of those who felt that somehow projective geometry was basic. It continued to look like the right approach when Lie's ideas appeared and there seemed to be a good connection between Lie's classification and the types of geometry organized by Klein.

Mathematicians could now ask why they had believed Euclidean geometry to be the only one when, in fact, many different geometries existed. The first to take up this question successfully was the German mathematician Moritz Pasch, who argued in 1882 that the mistake had been to rely too heavily on physical intuition. In his view an argument in mathematics should depend for its validity not on the physical interpretation of the terms involved but upon purely formal criteria. Indeed, the principle of duality did violence to the sense of geometry as a formalization of what one believed about (physical) points and lines; one did not believe that these terms were interchangeable.

The ideas of Pasch caught the attention of the German mathematician David Hilbert, who, with the French mathematician Henri Poincaré, came to dominate mathematics at the turn of the century. In wondering why it was that mathematics—and in particular geometry—produced correct results, he came to feel increasingly that it was not because of the lucidity of its definitions. Rather, mathematics worked because its (elementary) terms were meaningless. What kept it heading in the right direction was its rules of inference. Proofs were valid because they were constructed through the application of the rules of inference, according to which new assertions could be declared to be true simply because they could be derived, by means of these rules, from the axioms or previously proved theorems. The theorems and axioms were viewed as formal statements that expressed the relationships between these terms.

The rules governing the use of mathematical terms were arbitrary, Hilbert argued, and each mathematician could choose them at will, provided only that the choices made were self-consistent. A mathematician produced abstract systems unconstrained by the needs of science, and, if scientists found an abstract system that fit one of their concerns, they could apply the system secure in the knowledge that it was logically consistent.

Hilbert first became excited about this point of view (presented in his *Grundlagen der Geometrie* ["Foundations of Geometry"], 1899) when he saw that it led not merely to a clear way of sorting out the geometries in Klein's hierarchy according to the different axiom systems they obeyed but to new geometries as well. For the first time there was a way of discussing geometry that lay beyond even the very general terms proposed by Riemann. Not all of these geometries have continued to be of interest, but the general moral that Hilbert first drew for geometry he was shortly to draw for the whole of mathematics.

The foundations of mathematics. By the late 19th century the debates about the foundations of geometry had become the focus for a running debate about the nature of the branches of mathematics. Cauchy's work on the foundations of the calculus, completed by the German mathematician Karl Weierstrass in the late 1870s, left an edifice that rested on concepts such as that of the natural number (the integers 1, 2, 3, and so on) and on certain constructions involving them. The algebraic theory of numbers and the transformed theory of equations had focused attention on abstract structures in mathematics. Questions that had been raised about numbers since Babylonian times turned out to be best cast theoretically in terms of entirely modern creations whose independence from the physical world was beyond dispute. Finally, geometry, far from being a kind of abstract physics, was now seen as dealing with meaningless terms obeying arbitrary systems of rules. Although there had been no conscious

The axiomatization of geometry

The Erlanger Programm

plan leading in that direction, the stage was set for a consideration of questions about the fundamental nature of mathematics.

Similar currents were at work in the study of logic, which had also enjoyed a revival during the 19th century. The work of the English mathematician George Boole and the American Charles Sanders Peirce had contributed to the development of a symbolism adequate to explore all elementary logical deductions. Significantly, Boole's book on the subject was called *An Investigation of the Laws of Thought, on Which Are Founded the Mathematical Theories of Logic and Probabilities* (1854). In Germany, the logician Gottlob Frege had directed keen attention to such fundamental questions as what it means to define something and what sorts of purported definitions actually do define.

Cantor. All of these debates came together through the pioneering work of the German mathematician Georg Cantor on the concept of a set. Cantor had begun work in this area because of his interest in Riemann's theory of trigonometric series, but the problem of what characterized the set of all real numbers came to occupy him more and more. He began to discover unexpected properties of sets. For example, he could show that the set of all algebraic numbers, and a fortiori the set of all rational numbers, is countable in the sense that there is a one-to-one correspondence between the integers and the members of each of these sets by means of which for any member of the set of algebraic numbers (or rationals), no matter how large, there is always a unique integer it may be placed in correspondence with. But, more surprisingly, he could also show that the set of all real numbers is not countable. So, although the set of all integers and the set of all real numbers are both infinite, the set of all real numbers is a strictly larger infinity. This was in complete contrast to the prevailing orthodoxy, which proclaimed that infinite could only mean "larger than any finite amount."

Here the concept of number was being extended and undermined at the same time. The concept was extended because it was now possible to count and order sets that the set of integers was too small to measure; and it was undermined because even the integers ceased to be basic undefined objects. Cantor himself had given a way of defining real numbers as certain infinite sets of rational numbers. Rational numbers were easy to define in terms of the integers, but now integers could be defined by means of sets. One way was given by Frege in *Die Grundlagen der Arithmetik* (1884). He regarded two sets as the same if they contained the same elements. So in his opinion there was only one empty set (today symbolized by \emptyset), the set with no members. A second set could be defined as having only one element by letting that element be the empty set itself (symbolized by $\{\emptyset\}$), a set with two elements by letting them be the two sets just defined (*i.e.*, $\{\emptyset, \{\emptyset\}\}$), and so on. Having thus defined the integers in terms of the primitive concepts "set" and "element of," Frege agreed with Cantor that there was no logical reason to stop, and he went on to define infinite sets in the same way Cantor had. Indeed, Frege was clearer than Cantor about what sets and their elements actually were.

Frege's proposals went in the direction of a reduction of all mathematics to logic. He hoped that every mathematical term could be defined precisely and manipulated according to agreed, logical rules of inference. This, the "logician" program, was dealt an unexpected blow by the English mathematician and philosopher Bertrand Russell in 1902, who pointed out unexpected complications with the naive concept of a set. Nothing seemed to preclude the possibility that some sets were elements of themselves while others were not, but asked Russell, "what then of the set of all sets that were not elements of themselves?" If it is an element of itself, then it is not (an element of itself), but, if it is not, then it is—a paradox. Either the idea of a set as an arbitrary collection of already defined objects was flawed, or else the idea that one could legitimately form the set of all sets of a given kind was incorrect. Frege's program never recovered from this blow, and the theories of Russell, which he developed together with Alfred North Whitehead in their *Principia Mathematica* (1910–13), that

went in the same direction never found lasting appeal with mathematicians.

Greater interest attached to the ideas that Hilbert and his school began to advance. It seemed to them that what had worked once for geometry could work again for all of mathematics. Rather than attempt to define things so that problems could not arise, they suggested that it was possible to dispense with definitions and cast all of mathematics in an axiomatic structure using the ideas of set theory. Indeed, the hope was that the study of logic could be embraced in this spirit, thus making logic a branch of mathematics, the opposite of Frege's intention. There was considerable progress in this direction, and there emerged both a powerful school of mathematical logicians (notably in Poland) and an axiomatic theory of sets that avoided Russell's paradoxes and the others which had sprung up.

In the 1920s Hilbert put forward his most detailed proposal for establishing the validity of mathematics. According to his theory of proofs, everything was to be put into an axiomatic form, allowing the rules of inference to be only those of elementary logic, and only those conclusions that could be reached from this finite set of axioms and rules of inference were to be admitted. He proposed that a satisfactory system would be one which was consistent, complete, and decidable. By consistent Hilbert meant that it should be impossible to derive both a statement and its negation; by complete, that every properly written statement should be such that either it or its negation was derivable from the axioms; by decidable, that one should have an algorithm which determines of any given statement whether it or its negation is provable. Such systems did exist, for example, the first-order predicate calculus, but none had been found capable of allowing mathematicians to do interesting mathematics.

Hilbert's program, however, did not last long. In 1931 the Austrian-born American mathematician and logician Kurt Gödel showed that there was no system of Hilbert's type within which the integers could be defined and which was both consistent and complete. Later Gödel and, independently, the English mathematician Alan Turing showed that decidability was also unattainable. Perhaps paradoxically, the effect of this dramatic discovery was to alienate mathematicians from the whole debate. Instead, mathematicians, who may not have been too unhappy with the idea that there is no way of deciding the truth of a proposition automatically, learned to live with the idea that not even mathematics rests on rigorous foundations. Progress since has been in other directions. An alternative axiom system for set theory was later put forward by the Hungarian-born American mathematician John von Neumann, which he hoped would help resolve contemporary problems in quantum mechanics. There was also a renewal of interest in statements that are both interesting mathematically and independent of the axiom system in use. The first of these was the American mathematician Paul Cohen's surprising resolution in 1963 of the continuum hypothesis, which was Cantor's conjecture that the set of all subsets of the rational numbers was of the same size as the set of all real numbers. This turns out to be independent of the usual axioms for set theory, so there are set theories (and therefore types of mathematics) in which it is true and others in which it is false.

Mathematical physics. At the same time that mathematicians were attempting to put their own house in order, they were also looking with renewed interest at contemporary work in physics. The man who did the most to rekindle their interest was Poincaré. Poincaré showed that dynamic systems described by quite simple differential equations, such as the solar system, can nonetheless yield the most random-looking, chaotic behaviour. He went on to explore ways in which mathematicians can nonetheless say things about this chaotic behaviour and so pioneered the way in which probabilistic statements about dynamic systems can be found to describe what otherwise defies our intelligence.

Poincaré later turned to problems of electrodynamics. After many years' work, the Dutch physicist Hendrik Antoon Lorentz had been led to an apparent dependence of length and time on motion, and Poincaré was pleased to notice

The
axiomatic
school

The
continuum
hypothesis

Set theory

that the transformations that Lorentz proposed as a way of converting one observer's data into another's formed a group. This appealed to Poincaré and strengthened his belief that there was no sense in a concept of absolute motion; all motion was relative. Poincaré thereupon gave an elegant mathematical formulation of Lorentz's ideas, which fitted them into a theory in which the motion of the electron is governed by Maxwell's equations. Poincaré, however, stopped short of denying the reality of the ether or of proclaiming that the velocity of light is the same for all observers, so credit for the first truly relativistic theory of the motion of the electron rests with Einstein and his special theory of relativity (1905).

Einstein's special theory is so called because it treats only the special case of uniform relative motion. The much more important case of accelerated motion and motion in a gravitational field was to take a further decade and to require a far more substantial dose of mathematics. Einstein only changed his estimate of the value of pure mathematics, which he had hitherto disdained, when he discovered that many of the questions he was led to had already been formulated mathematically and had been solved. He was most struck by theories derived from the study of geometry in the sense in which Riemann had formulated it.

By 1915 a number of mathematicians were interested in reapplying their discoveries to physics. The leading institution in this respect was the University of Göttingen, where Hilbert had unsuccessfully attempted to produce a general theory of relativity before Einstein, and it was there that many of the leaders of the coming revolution in quantum mechanics were to study. There, too, went many of the leading mathematicians of their generation, notably John von Neumann and Hermann Weyl, to study with Hilbert. In 1904 Hilbert had turned to the study of integral equations. These arise in many problems where the unknown is itself a function of some variable, and especially in those parts of physics that are expressed in terms of extremal principles (such as the principle of least action). The extremal principle usually yields information about an integral involving the sought-for function, hence the name "integral equation." Hilbert's contribution was to bring together many different strands of contemporary work and to show how they could be elucidated if cast in the form of arguments about objects in certain infinite-dimensional vector spaces.

The extension to infinite dimensions was not a trivial task, but it brought with it the opportunity to use geometric intuition and geometric concepts to analyze problems about integral equations. Hilbert left it to his students to provide the best abstract setting for his work, and thus was born the concept of a Hilbert space. Roughly speaking this is an infinite-dimensional vector space in which it makes sense to speak of the lengths of vectors and the angles between them; useful examples include certain spaces of sequences and certain spaces of functions. Operators defined on these spaces are also of great interest: their study forms part of the field of functional analysis.

When in the 1920s mathematicians and physicists were seeking ways to formulate the new quantum mechanics, von Neumann proposed that the subject be written in the language of functional analysis. The quantum mechanical world of states and observables, with its mysterious wave packets that were sometimes like particles and sometimes like waves depending on how they were observed, went very neatly into the theory of Hilbert spaces. Functional analysis has ever since grown with the fortunes of particle physics.

Algebraic topology. The early 20th century saw the emergence of a number of theories whose power and utility reside in large part in their generality. Typically, they are marked by an attention to the set or space of all examples of a particular kind. (Functional analysis is such an endeavour.) One of the most energetic of these general theories was that of algebraic topology. In this subject a variety of ways are developed for replacing a space by a group and a map between spaces by a map between groups. It is like using X rays; information is lost, but the shadowy image of the original space may turn out

to contain, in an accessible form, enough information to solve the question at hand.

Interest in this kind of research came from various directions. Galois's theory of equations was an example of what could be achieved by transforming a problem in one branch of mathematics into a problem in another, more abstract branch. Another impetus came from Riemann's theory of complex functions. He had studied algebraic functions, that is, loci defined by equations of the form $f(x, y) = 0$, where f is a polynomial in x whose coefficients are polynomials in y . When x and y are complex variables, the locus can be thought of as a real surface spread out over the x plane of complex numbers (today called a Riemann surface). To each value of x there correspond a finite number of values of y . Such surfaces are not easy to comprehend, and Riemann had proposed to draw curves along them in such a way that, if the surface was cut open along them, it could be opened out into a polygonal disk (see Figure 18). He was able to establish a profound connection between the minimum number of curves needed to do this for a given surface and the number of functions (becoming infinite at specified points) that the surface could then support.

The natural problem was to see how far Riemann's ideas could be applied to the study of spaces of higher dimension. Here two lines of inquiry developed. One emphasized what could be obtained from looking at the projective geometry involved. This point of view was fruitfully applied by the Italian school of algebraic geometers. It ran into problems, which it was not wholly able to solve, having to do with the singularities a surface can possess. Whereas a locus given by $f(x, y) = 0$ may intersect itself only at isolated points, a locus given by an equation of the form $f(x, y, z) = 0$ may intersect itself along curves, a problem that caused considerable difficulties. The second approach emphasized what can be learned from the study of integrals along paths on the surface. This approach, pursued by Charles-Émile Picard and by Poincaré, provided a rich generalization of Riemann's original ideas.

On this base, conjectures were made and a general theory produced, first by Poincaré and then by the American engineer-turned-mathematician Solomon Lefschetz (1884–1972), concerning the nature of manifolds of arbitrary dimension. Roughly speaking a manifold is the n -dimensional generalization of the idea of a surface: it is a space any small piece of which looks like a piece of n -dimensional space. Such an object is often given by a single algebraic equation in $n + 1$ variables. At first their work was concerned with how these manifolds may be decomposed into pieces, counting the number of pieces and decomposing them in their turn. The result was a list of numbers, called Betti numbers in honour of the Italian mathematician Enrico Betti, who had taken the first steps of this kind to extend Riemann's work. It was only in the late 1920s that the German mathematician Emmy Noether suggested how the Betti numbers might be thought of as measuring the size of certain groups. At her

Riemann surfaces

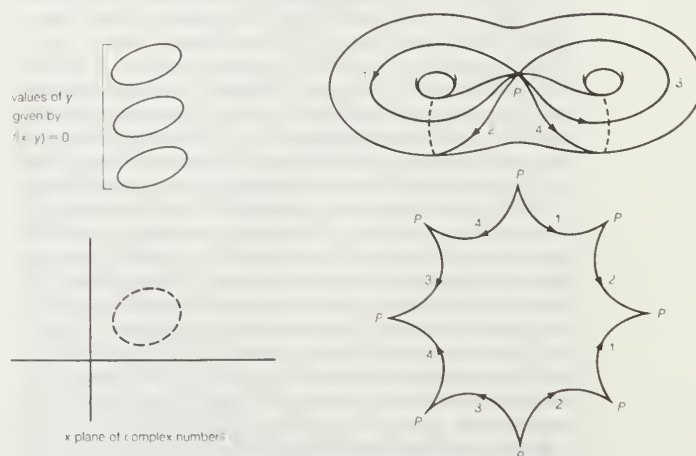


Figure 18: (Left) Pieces of a surface given by $f(x, y) = 0$; (right) if the surface is cut along the curves, an octagon is obtained

Integral equations

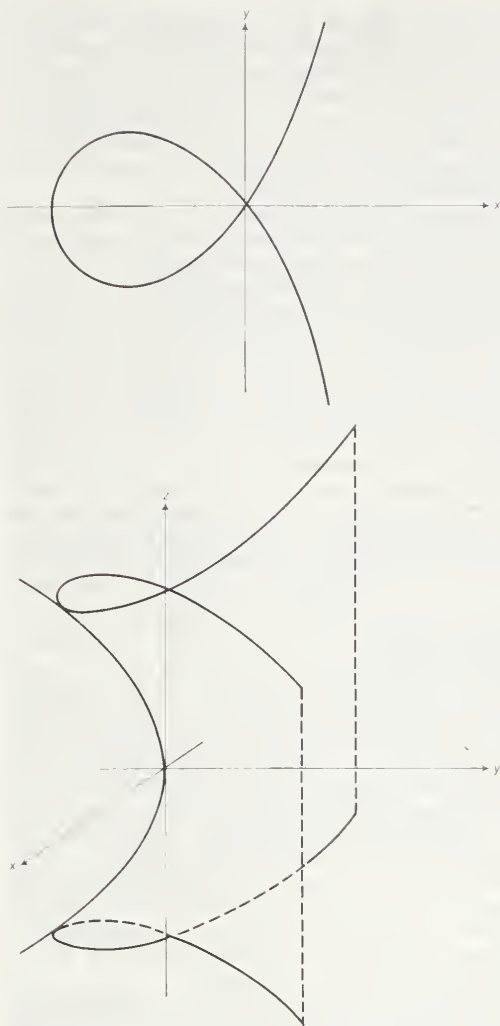


Figure 19: (Top) $f(x, y) = x^2(x + 1) - y^2 = 0$ intersects itself at $(x, y) = (0, 0)$; (bottom) $E(x, y, z) = 0 = y^2(y + z^2) - x^2$ intersects itself along the z -axis, but the origin is a triple self-intersection.

instigation a number of people then produced a theory of these groups, the so-called homology and cohomology groups of a space.

Homology and cohomology

Two objects that can be deformed into one another will have the same homology and cohomology groups. To assess how much information is lost when a space is replaced by its algebraic topological picture, Poincaré asked the crucial converse question: according to what algebraic conditions is it possible to say that a space is topologically equivalent to a sphere? He showed by an ingenious example that having the same homology is not enough and proposed a more delicate index, which has since grown into the branch of topology called homotopy theory. Being more delicate, it is both more basic and more difficult. There are usually standard methods for computing homology and cohomology groups, and they are completely known for many spaces. In contrast, there is scarcely an interesting class of spaces for which all the homotopy groups are known. And Poincaré's original question continues to resist answer, but only in the dimension in which he raised it: there is still no algebraic criterion for recognizing the 3-sphere.

Developments in pure mathematics. The interest in axiomatic systems at the turn of the century led to axiom systems for the known algebraic structures, that for the theory of fields, for example, being developed by the German mathematician Ernst Steinitz in 1910. The theory of rings (structures in which it is possible to add, subtract, and multiply but not necessarily divide) was much harder to formalize. It is important for two reasons: the theory of algebraic integers forms part of it, because algebraic integers naturally form into rings; and (as Kronecker and Hilbert had argued) algebraic geometry forms another

part. The rings that arise there are rings of functions definable on the curve, surface, or manifold or are definable on specific pieces of it.

Problems in number theory and algebraic geometry are often very difficult, and it was the hope of mathematicians such as Noether, who laboured to produce a formal, axiomatic theory of rings, that by working at a more rarefied level the essence of the concrete problems would remain, while the distracting special features of any given case would fall away. This would make the formal theory both more general and easier, and to a surprising extent these mathematicians were successful.

A further twist to the development came with the work of the American mathematician Oscar Zariski (1899–1986), who had studied with the Italian school of algebraic geometers but came to feel that their method of working was imprecise. He worked out a detailed program whereby every kind of geometric configuration could be redescribed in algebraic terms. His work succeeded in producing a rigorous theory, although some, notably Lefschetz, felt that the geometry had been lost sight of in the process.

The study of algebraic geometry was amenable to the topological methods of Poincaré and Lefschetz so long as the manifolds were defined by equations whose coefficients were complex numbers. But with the creation of an abstract theory of fields it was natural to want a theory of varieties defined by equations with coefficients in an arbitrary field. This was provided for the first time by the French mathematician André Weil (b. 1906), in his *Foundations of Algebraic Geometry* (1946), in a way that drew on Zariski's work without suppressing the intuitive appeal of geometric concepts. Weil's theory of polynomial equations is the proper setting for any investigation that seeks to determine what properties of a geometric object can be derived solely by algebraic means. But it falls tantalizingly short of one topic of importance: the solution of polynomial equations in integers. This was the topic that Weil took up next.

Weil's theory of polynomial equations

The central difficulty is that in a field it is possible to divide, but in a ring it is not. The integers form a ring but not a field (dividing 1 by 2 does not yield an integer). But Weil showed that simplified versions (posed over a field) of any question about integer solutions to polynomials could be profitably asked. This transferred the questions to the domain of algebraic geometry. To count the number of solutions, Weil proposed that, since the questions were now geometric, they should be amenable to the techniques of algebraic topology. This was an audacious move, since there was no suitable theory of algebraic topology available, but Weil conjectured what results it should yield. The difficulty of Weil's conjectures may be judged by the fact that the last of them was a generalization to this setting of the famous Riemann hypothesis about the zeta function, and they rapidly became the focus of international attention.

Weil, along with Claude Chevalley, Henri Cartan, Jean Dieudonné, and others, created a group of young French mathematicians who began to publish virtually an encyclopaedia of mathematics under the name Nicolas Bourbaki, taken by Weil from an obscure general of the Franco-German War. Bourbaki became a self-selecting group of young mathematicians who were strong on algebra, and the individual Bourbaki members were interested in the Weil conjectures. In the end, they succeeded completely. A new kind of algebraic topology was developed, and the Weil conjectures were proved. The generalized Riemann hypothesis was the last to surrender, being established by the Belgian mathematician Pierre Deligne in the early 1970s. Strangely, its resolution still leaves the original Riemann hypothesis unsolved.

Bourbaki was a key figure in the rethinking of structural mathematics. Algebraic topology was axiomatized by Samuel Eilenberg, a Polish-born American mathematician and Bourbaki member, and the American mathematician Norman Steenrod. Saunders MacLane, also of the United States, and Eilenberg extended this axiomatic approach until many types of mathematical structures were presented in families, called categories. Hence there was a category consisting of all groups and all maps between

Categories

them that preserve multiplication, and there was another category of all topological spaces and all continuous maps between them. To do algebraic topology was to transfer a problem posed in one category (that of topological spaces) to another (usually that of commutative groups or rings). When he created the right algebraic topology for the Weil conjectures, the German-born French mathematician Alexandre Grothendieck, a Bourbaki of enormous energy, produced a new description of algebraic geometry. In his hands it became infused with the language of category theory. The route to algebraic geometry became steeper than ever, but the views from the summit have a naturalness and a profundity that have brought many experts to prefer it to the earlier formulations, including Weil's.

Grothendieck's formulation makes algebraic geometry the study of equations defined over rings rather than fields. Accordingly, it raises the possibility that questions about the integers can be answered directly. Building on the work of like-minded mathematicians in the United States, France, and Russia, the German Gerd Faltings triumphantly vindicated this approach when he solved the English mathematician Louis Mordell's conjecture in 1983. This conjecture states that almost all polynomial equations that define curves have at most finitely many rational solutions; the cases excluded from the conjecture are the simple ones that are much better understood.

Meanwhile, the German mathematician Gerhard Frey had pointed out that, if Fermat's last theorem is false, so that there are integers u , v , w such that $u^p + v^p = w^p$ (p greater than 5), then for these values of u , v , and p the curve $y^2 = x(x-u^p)(x+v^p)$ has properties that contradict major conjectures of the Japanese mathematicians Taniyama Yutaka and Shimura Goro about elliptic curves. Frey's observation, refined by the French mathematician Jean-Pierre Serre and proved by the American mathematician Ken Ribet, meant that by 1990 Taniyama's unproved conjectures were known to imply Fermat's last theorem.

In 1993 the English mathematician Andrew Wiles established the Shimura-Taniyama conjectures in a large range of cases that included Frey's curve and therefore Fermat's last theorem—a major feat even without the connection to Fermat. It soon became clear that the argument had a serious flaw; but in May 1995 Wiles, assisted by another English mathematician, Richard Taylor, published a different and valid approach. In so doing, Wiles not only solved the most famous outstanding conjecture in mathematics but also triumphantly vindicated the sophisticated and difficult methods of modern number theory.

Mathematical physics and the theory of groups. In the 1910s the ideas of Lie and Killing were taken up by the French mathematician Élie-Joseph Cartan, who simplified their theory and rederived the classification of what came to be called the classical complex Lie algebras. The simple Lie algebras, out of which all the others in the classification are made, were all representable as algebras of matrices, and in a sense Lie algebra is the abstract setting for matrix algebra. Connected to each Lie algebra there were a small number of Lie groups, and there was a canonical simplest one to choose in each case. The groups had an even simpler geometric interpretation than the corresponding algebras, for they turned out to describe motions that leave certain properties of figures unaltered. For example, in Euclidean three-dimensional space, rotations leave unaltered the distances between points; the set of all rotations about a fixed point turns out to form a Lie group, and it is one of the Lie groups in the classification. The theory of Lie algebras and Lie groups shows that there are only a few sensible ways to measure properties of figures in a linear space and that these methods yield groups of motions leaving the figures, which are (more or less) groups of matrices, unaltered. The result is a powerful theory that could be expected to apply to a wide range of problems in geometry and physics.

The leader in the endeavours to make Cartan's theory, which was confined to Lie algebras, yield results for a corresponding class of Lie groups was the German-American mathematician Hermann Weyl. He produced a rich and satisfying theory for the pure mathematician and wrote extensively on differential geometry and group theory and

its applications to physics. Weyl attempted to produce a theory that would unify gravitation and electromagnetism. His theory met with criticism from Einstein and was generally regarded as unsuccessful; it has been only in the last quarter of the 20th century that similar unified field theories have met with any acceptance. Nonetheless, Weyl's approach demonstrates how the theory of Lie groups can enter into physics in a substantial way.

In any physical theory the endeavour is to make sense of observations. Different observers make different observations. If they differ in choice and direction of their coordinate axes, they give different coordinates to the same points, and so on. Yet the observers agree on certain consequences of their observations: in Newtonian physics and Euclidean geometry they agree on the distance between points. Special relativity explains how observers in a state of uniform relative motion differ about lengths and times but agree on a quantity called the interval. In each case they are able to do so because the relevant theory presents them with a group of transformations that converts one observer's measurements into another's and leaves the appropriate basic quantities invariant. What Weyl proposed was a group that would permit observers in nonuniform relative motion, and whose measurements of the same moving electron would differ, to convert their measurements and thus permit the (general) relativistic study of moving electric charges.

In the 1950s the American physicists Chen Ning Yang and Robert L. Mills gave a successful treatment of the so-called strong interaction in particle physics from the Lie group point of view. Twenty years later mathematicians took up their work, and a dramatic resurgence of interest in Weyl's theory began. These new developments, which had the incidental effect of enabling mathematicians to escape the problems in Weyl's original approach, were the outcome of lines of research that had originally been conducted with little regard for physical questions. Not for the first time, mathematics was to prove surprisingly or, as the Hungarian-born American physicist Eugene Wigner said, "unreasonably effective" in science.

Cartan had investigated how much may be accomplished in differential geometry using the idea of moving frames of reference. This work, which was partly inspired by Einstein's theory of general relativity, was also a development of the ideas of Riemannian geometry that had originally so excited Einstein. In the modern theory one imagines a space (usually a manifold) made up of overlapping coordinatized pieces. On each piece, one supposes some functions to be defined, which might in applications be the values of certain physical quantities. Rules are given for interpreting these quantities where the pieces overlap. The data are thought of as a bundle of information provided at each point. For each function defined on each patch, it is supposed that at each point a vector space is available as mathematical storage space for all its possible values. Because a vector space is attached at each point, the theory is called the theory of vector bundles. Other kinds of space may be attached, thus entering the more general theory of fibre bundles. The subtle and vital point is that it is possible to create quite different bundles which nonetheless look similar in small patches. An example of this is illustrated in Figure 20. The cylinder and the Möbius band look alike in small pieces but are topologically distinct, since it is possible to give a standard sense of direction to all the lines in the cylinder but not to those in the Möbius band. Both spaces can be thought of as one-dimensional vector bundles over the circle, but they are very different. The cylinder is regarded as a "trivial" bundle, the Möbius band as a twisted one.

In the 1940s and '50s a vigorous branch of algebraic topology established the main features of the theory of bundles. Then, in the 1960s, work chiefly by Grothendieck and the English mathematician Michael Atiyah showed how the study of vector bundles on spaces could be regarded as the study of cohomology theory (called K theory). More significantly still, in the 1960s Atiyah, the American Isadore Singer, and others found ways of connecting this work to the study of a wide variety of questions involving partial differentiation, culminating in the celebrated

The theory of vector bundles

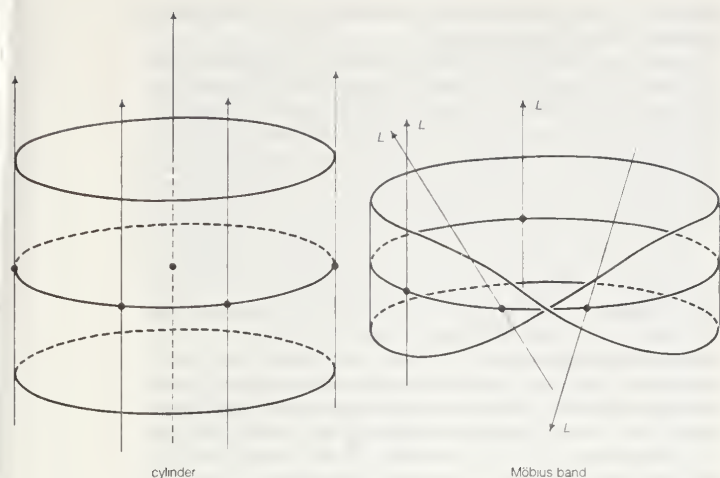


Figure 20: Vector bundles.

As the circle is followed clockwise around the Möbius band, the line L twists through the half a turn, so the lines cannot be consistently made to point in the same direction.

Atiyah–Singer theorem for elliptic operators. (Elliptic is a technical term for the type of operator studied in potential theory.) There are remarkable implications for the study of pure geometry, and much attention has been directed to the problem of how the theory of bundles embraces the theory of Yang and Mills, which it does precisely because there are nontrivial bundles, and to the question of how it can be made to pay off in large areas of theoretical physics. These include the theories of superspace and supergravity and the string theory of fundamental particles, which involves the theory of Riemann surfaces in novel and unexpected ways. (J.J.G.)

Mathematics in China and Japan

When speaking of mathematics in East Asia, it is necessary to take into account its development in China, Korea, and Japan as a whole. Mathematicians from these countries can be considered as part of the same working community. Moreover, these scholars usually wrote with Chinese characters and thus could read one another's texts.

The common basis for the development of mathematics in East Asia was books that were written in China from the 1st through the 13th century, to which most of the subsequent works refer. For this period, the two main sources, the references found in the texts and those given in the bibliographies compiled for the dynastic annals, indicate that there are many lacunae in the books that have survived. The oldest works extant survived because they became official books.

In the 17th century few ancient Chinese mathematical works were known, and those that were known were not fully understood. Thereafter, as Chinese mathematicians became aware of European achievements, they began to look for such works throughout the country. Editions of the texts they found began to appear at the end of the 18th century and have become the main sources for the history of Chinese mathematics. Discoveries of new sources are now rare, contrary to the situation for Arabic mathematics, for example. Nevertheless, in the 20th century a mathematical book was discovered in a grave dating from the 2nd or 3rd century, which pushed back by more than three centuries what was previously known about the subject.

CHINESE MATHEMATICS TO THE 13TH CENTURY

Outline of the history. The *Nine Chapters on the Mathematical Procedures* (or the *Nine Chapters*) was probably compiled in the 1st century AD. This book gathered and organized many mathematical achievements from preceding periods. It played an important part in the development of mathematics in China, for all Chinese mathematicians refer to it and most of the subjects they have worked on stem from it. Its format, which was adopted by most subsequent authors, consists of problems for which a numerical

answer and a procedure for solution are given. These problems are arranged in classes that come under the heading of a general method. Many scholars wrote commentaries on the *Nine Chapters*, adding explanations and proofs, rewriting procedures, and suggesting new formulas. The most important of the commentaries to survive, attributed to Liu Hui (3rd century), contains the richest set of proofs within this tradition.

Some of the books written subsequently are known because, gathered together with the *Nine Chapters* and the astronomical treatise *Mathematical Classic of the Gnomon of Chou* and commented on in the 7th century by a group under the leadership of Li Ch'un-feng, they became the *Ten Classics of Mathematics*, the manual for officials trained in the then newly established office of mathematics. Although some people were thus officially trained as mathematicians, no major breakthrough seems to have been achieved until the 11th century.

At that time (1084) the *Ten Classics of Mathematics* was edited and printed, and this seems to have been related to renewed activity in mathematics during the 11th and 12th centuries, known today only through later quotations. This period probably paved the way for the major achievements of Chinese tradition, as they are known today only through the few books that have come down from the second half of the 13th century. China was then divided into two countries, North and South, and achievements by mathematicians of both sides are known: in the South, those of Ch'in Chiu-shao and Yang Hui (who were minor officials), and in the North, those of Li Yeh (a recluse scholar) and Chu Shih-chieh (a wandering teacher). Their contributions seem to have been arrived at independently but they attest to a common basis.

From the period after this, no valuable works survive, and there is no evidence that any important mathematical works were written. Still, some major works of the 13th century are recorded in the encyclopaedia compiled under the Yung-lo emperor (1402–24), but commentaries on these books written by the end of the 15th century show that by this time they were no longer understood.

It was only in the 16th century that the abacus appears to have come into widespread use, and most of the books of the period discuss it. But it is not possible to date the moment when it actually appeared in China.

The following discussion of the evolution of mathematical subjects within the Chinese tradition emphasizes several common characteristics of most of the achievements: a specific use of algorithms and the importance given to position, to configurations of numbers, and to parallelisms between procedures.

The Nine Chapters. This book presupposes mathematical knowledge about how to represent numbers and how to perform the four arithmetic operations. The numbers are written in Chinese characters, but, for most of the procedures described, the computations are to be performed on a surface, perhaps on the ground (see below). Most probably, as can be inferred from later accounts, the numbers were represented with counting rods, according to a decimal place-value system. Such a representation can be moved and modified within a computation. Indian and Arabic mathematicians computed (among other ways) on erasable surfaces, and it appears that arithmetic operations were performed in China in ways comparable to these, although no written computations were recorded until much later. As will be seen, setting up the computations in this way greatly influenced later mathematical developments.

The *Nine Chapters* contains a number of mathematical achievements, already in a mature form, that are presented by most of the subsequent books without substantial changes.

Arithmetic of fractions. Division is a central operation in the *Nine Chapters*. The numbers used throughout are of the type that results from division: "integer plus fraction." Fractions are defined to be a natural part of the result of a division, the remainder of the division being taken as the numerator and the divisor as the denominator. Thus, dividing 17 by 5, one obtains 3 as the quotient and a remainder of 2. This gives rise to the fraction $\frac{2}{5}$. The fractions are thus always less than one, and their arithmetic is

Computational methods

described through the use of division. To get the sum of a set of fractions, one is instructed to "multiply the numerators by the denominators that do not correspond to them, add to get the dividend. Multiply the denominators all together to get the divisor. Perform the division. If there is a remainder, name it with the divisor." The sum of a set of fractions is itself thus the result of a division, of the form integer plus fraction. All the operations are described in a similar way. The book also contains explanations of such algorithms as the "rule of three," "sharing in unequal parts," "false double position," and so on.

Formulas for the computations of areas and volumes. The *Nine Chapters* gives formulas for elementary plane and solid figures as well as for the area of the circle, segment of a circle, and sphere. All these formulas are expressed by lists of operations to perform on the data to get the result, an algorithmic expression that is common in ancient traditions (see below on the remarkable elaboration of algorithms found in the book). For example, to compute the area of the circle, the following algorithm is given: "multiply the diameter by itself, triple this, divide by four." This algorithm amounts to taking 3 for π . Commentators added improved values for π along with some derivations. Liu Hui gave two other values for π , one slightly low ($157/50$) and one high ($3927/1250$). The *Nine Chapters* also provides the correct formula for the area of the circle: "multiplying half the diameter and half the circumference, one gets the area," which Liu Hui proved.

Solution of systems of simultaneous linear equations. The chapter devoted to this topic provides a general procedure. For instance, one of its problems, on the yields from three grades of grain, leads to a system of three linear equations in three unknowns, and the procedure for the solution arranges the data on the counting surface in a table, as in Figure 21. The coefficients of the first equation are arranged in the right column. The coefficients of the second and third equations are arranged in the middle and left columns, respectively, and so on. Consequently, the numbers of the first row all correspond to the first unknown. This is an instance of a positional notation, in which the position of a number in a numerical configuration has a mathematical meaning. The main tool for the solution is the use of column reduction to obtain an equivalent configuration (see Figure 21). Next, the unknown of the third line is found by division, and hence the second and the first unknowns as well.

The description of the algorithm uses the structure that the configuration gives to the data on the counting surface in an essential way. Because this procedure implies a column-to-column subtraction, it gives rise to negative numbers. Detailed methods of computing with positive and negative coefficients are used to solve problems involving two to seven unknowns. This seems to be the first occurrence of negative numbers in the history of mathe-

atics. The *Nine Chapters* also contains some knowledge that was to be reworked later.

Square root and cube root extraction. Algorithms in the *Nine Chapters* for finding integer parts of roots on the counting surface are based on the same idea as the arithmetic ones used today. They are described using the technical terms of division, in such a way that both roots appear to be derived from this operation. These algorithms are set up on the surface in the same way as is a division ("quotient" on the top; under it, the "dividend"; one row below, the "divisor"; at the bottom, auxiliary computations). Moreover, the two algorithms are written out, sentence by sentence, parallel to each other, so that their similarities and differences become clear.

In the writing of algorithms, there are three basic operations: iteration, conditional statements, and assignation of variables. Once the integer parts are found, the algorithms supply a fractional approximation of the root. (To find the square root of A, for example, if the integer part of the root found by the algorithm is a, the result is given as $[a + (A - a^2)/2a]$.)

In commenting on these algorithms, Liu Hui gives new approximations and suggests that one should continue computing the digits in the same way, setting 10 as denominator for the first digit, 100 as denominator for the second digit, and so on; he thus gives the root in terms of decimal fractions.

In the *Nine Chapters* a problem is solved with a quadratic equation, which appears to be conceived of as an arithmetic operation depending on square root extraction. The procedure for solving the equation is part of the procedure for root extraction. However, the equation is thought to have only one root. The theory of equations developed in China within that framework, until its apogee in the 13th century. (The solution by radicals that Babylonian mathematicians had already explored has not been found in the Chinese texts that survive.)

Problems involving right-angled triangles. The so-called Pythagorean theorem is given in the *Nine Chapters*. Formulas are provided to solve problems on right-angled triangles such as the following: "given the base, and the sum of the height and of the hypotenuse, find the height and the hypotenuse." The diameter of the inscribed circle and the side of the inscribed square are given as well.

The Commentary of Liu Hui. Liu Hui's commentary on the *Nine Chapters* is the most important text dating from before the 13th century that contains proofs. He gives proofs for the formulas of volumes that are presented in the *Nine Chapters*, and he adds new formulas for the same volumes. He also organizes these formulas, given one after the other without comment in the *Nine Chapters*, into a system in which proofs for one formula use only formulas that have been established independently. Starting from the formula of the parallelepiped, Liu Hui proves the formulas of "standard blocks," from which he then deduces other formulas. He proceeds with a small set of proof techniques, including dissection, decomposition and recomposition into known pieces, and Cavalieri's principle (which states that, if two solids of the same height are such that their corresponding sections at any level have the same areas, then they have the same volume), and so on.

Liu Hui compares the algorithms given in the algebraic sections of the *Nine Chapters* with one another and demonstrates how the same formal operations, which he calls the "key-steps" of computations, are brought into play in different algorithms. Thus, in the comparison between the addition of fractions and the solution of systems of simultaneous linear equations mentioned above, Liu Hui shows that sets of numbers are involved (numerator and denominator for a fraction, the coefficients of an equation for systems of equations) which share the property that all the numbers of a set can be multiplied by the same number without altering the mathematical meaning of that set. Both algorithms proceed by multiplying the sets of numbers that enter into a problem, each by an appropriate factor, in such a way that some corresponding numbers of the sets are made equal and the other numbers multiplied to keep intact the meaning of the whole

Positional notation

1	2	3	$x + 2y + 3z = 26$
2	3	2	represents $2x + 3y + z = 34$
3	1	1	$3x + 2y + z = 39$
26	34	39	

Proof techniques

1	6	3	1	3	3
2	9	2	2	7	2
3	3	1	3	2	1
26	34	39	26	63	39

Figure 21: The first example of a system of linear equations in the *Nine Chapters*.

(Top) The coefficients of each equation are arranged in columns in the box; the coefficients of the middle column are multiplied by 3 (the first position in the right column), and the right column is then subtracted as many times as necessary to produce 0 in the first position of the middle column. (Bottom) The operation is repeated to obtain a "triangular matrix" from which the unknowns are then obtained one after the other.

sets. (In the case of fractions, the denominators are made equal and the numerators changed appropriately; for systems of equations, the procedure is the same as if two numbers in the same row but in different columns were made equal by an appropriate multiplication, so that one of them can be eliminated through a column-to-column subtraction; the whole columns are then multiplied by the same number so that the equations remain true). Liu Hui proceeds from these analogies to state new algorithms for the same problems.

The Ten Classics of Mathematics. The *Ten Classics of Mathematics* contains explications of the mathematical knowledge presupposed by the *Nine Chapters* (the numeration system, arithmetic operations, etc.). Most of the subjects presented rely on the algorithms presented in the *Nine Chapters*: the rule of three, sharing in unequal parts, false double position, systems of simultaneous linear equations, and so on. These algorithms are not stated in as full generality in the *Ten Classics* as they are in the *Nine Chapters*, but they are given in the solutions of particular problems. Nevertheless, it is possible to see the ongoing evolution of some of these subjects, for example, the root extraction, the solution of equations, and the summation of series. The mathematical classics by Sun Tzu and Chang Ch'iu-chien (probably written before the 5th century) employ new ways of describing the algorithms of square root and cube root extraction. The underlying procedures are the same and they are still described in parallel ways, but this parallelism shows more clearly than before the mathematical object under which both of them are actually subsumed and which is thus responsible for their similarity: namely, equation. What has changed in the description is that, just as division involves a single divisor, square root extraction is shown to have two divisors, and cube root extraction three divisors. (These divisors actually are coefficients of the equations underlying the root extractions.) The divisors are shown to play similar roles in the algorithms. Moreover, in setting up the algorithms, the divisors are arranged one above the other, yielding a positional notation for the underlying equations: the row in which a number occurs is associated with the power of the unknown whose coefficient it is.

The *Ten Classics* contains discussions of topics that are not mentioned in the *Nine Chapters* but that were to be the subject of some of the highest mathematical achievements of the Sung and Yüan dynasties (960–1368). For example, the *Mathematical Classics* by Sun Tzu presents this congruence problem: Suppose one has an unknown amount of objects. If one counts them by 3s, there remain 2 of them. If one counts them by 5s, there remain 3 of them. If one counts them by 7s, there remain 2 of them. How many objects are there? The procedure used to solve the problem is difficult to understand, however, because it is described in a very condensed manner.

Some major developments from the 11th century to the 13th century. *Theory of root extraction and of equations.* In the 11th century Chia Hsien is said to have given an algorithm to find a fourth root involving a method similar to the so-called Ruffini–Horner method. In this algorithm the successive numbers that are put in each of the rows (actually the coefficients) are obtained through computations that involve only the numbers located in the rows below. Again the algorithm makes use of the configuration given to this set of numbers in an essential way. In addition, the procedures used to compute the numbers of any row are basically the same. The representation of its computation still involves a positional notation for the underlying equations. This new algorithm can be applied for square and cube root extraction as well; it is only necessary to determine the number of rows (*i.e.*, coefficients) that each of these computations involves. It was again through a reformulation of the previous root extraction procedures that a new improvement was obtained: these procedures were shown to be special cases of the same general algorithm.

The 12th-century scholar Liu I explored a method of finding roots of quadratic equations that have positive or negative coefficients; such equations arise in geometric problems. These coefficients, whatever their sign, are en-

tered in the table of the root extraction, and the square root algorithm is adapted to each situation.

Later, in Ch'in Chiu-shao's *Mathematics in Nine Chapters* (1247), a similar algorithm is used to find "the" root of any equation. By that time, general equations were used and were represented by a positional notation (see the example in the next section), which seems to indicate that it was the slow evolution of the algorithms of root extraction and their comparison that produced the concept of equation. Similar methods (with a slightly different notation) were well known to Li Yeh, and in his *Sea-Mirror of Circle Measurements*, written only one year after Ch'in completed his book, he takes the search for the root of equations for granted. Li Yeh lived in North China, while Ch'in Chiu-shao lived in the South, and he is thought to have worked without knowing Ch'in's achievement. It is thus highly probable that these methods were well known before the middle of the 13th century.

From Chia Hsien on, another method was known for finding the root of an equation, using the coefficients of what is now called Pascal's triangle (see Figure 22) and the same positional representation.

By permission of the Syndics of Cambridge University Library



Figure 22: A Chinese representation of Pascal's triangle.

The method of the "celestial unknown." Li Yeh's book also contains a method, unknown to Ch'in Chiu-shao, which seems to have flourished in North China for some decades. This method explains how to use polynomial arithmetic to find equations to solve a problem. Li Yeh's book is the oldest surviving work that explains this method, but it was probably not the first to deal with it. In this book polynomials are also arranged according to a positional notation. Thus $x^2 - 3x + 5 + 7/x^2$ is represented as

$$\begin{array}{r} 7 \\ 0 \\ 5 \\ -3 \\ 1 \end{array}$$

A character is added next to the 5 to indicate that it is a constant term. The location of the coefficient indicates the power of the indeterminate with which it is associated. This indeterminate is called the celestial unknown.

It is known that some mathematicians used this representation for polynomials in two or three unknowns. In his book *The Jade Mirror of the Four Unknowns*, Chu Shih-chieh makes use of four unknowns. Starting from the centre, in the two horizontal and the two vertical directions, he puts in increasing order of their powers what comes from each of the four unknowns. In such problems, where there is more than one unknown, he has to use a method of elimination of a common unknown to two equations.

Indeterminate analysis. Ch'in Chiu-shao's book also contains algorithms for the general congruence problem, some examples of which are given in Sun Tzu's treatise, where its solution was too obscure to be understood. This

Poly-nomials in several unknowns

Congru-
ence
problems

problem amounts to determining a number, the remainders of which are known when it is divided by given numbers (called moduli). There is no extant work between Sun's treatise and Ch'in's book of 1247 that reveals how this algorithm was elaborated. Such problems seem to have been worked out because of calendrical computation. Ch'in Chiu-shao introduces his discussion of these problems by saying that his goal is to clarify several procedures used by astronomers who were applying them without understanding them. His solution is known today as the Chinese remainder theorem. He deals with the case when moduli are relatively prime, and then reduces the case when they are not to it. The first case is easily solved when x can be found that satisfies the congruence $xa \equiv 1 \pmod{b}$, a and b being two given relatively prime numbers (suppose $a < b$). Ch'in gives an algorithm for this, using a sequence of quotients in searching for the greatest common divisor of a and b , which is also the sequence of convergents for the continued fraction for b/a . Having them, he is then able to compute x .

In addition, the summation of series was developed to a greater extent during this period.

The decline of the Sung-Yüan mathematics. Little is known about what happened after Chu Shih-chieh. In the 16th century a mathematician commenting on Li Yeh's *Sea-Mirror of Circle Measurements* no longer understood the method of the celestial unknown. By the 17th century it seems to have been completely forgotten. Rods were then no longer used as a counting tool, so perhaps Chinese algebraic positional notations, deprived of the instrument on which they were based, could not be understood. On the other hand, there was a rapid diffusion of the abacus, for which many books were written, including the *Systematic Treatise on Arithmetic* by Ch'eng Ta-wei (1592). This book explains in detail arithmetic on the abacus. Editions of this work remained popular until the 20th century.

When the Jesuits arrived in China at the end of the 16th century, they found people interested in science (so that they were accepted in China because of their scientific knowledge) but unaware of the Chinese past in mathematics. An era of translations of Western works started then (the six first books of Euclid's *Elements* were translated by a Jesuit and a Chinese in 1607), followed by a period when Chinese mathematicians attempted to find ancient books, to understand them, and to make a synthesis of the Chinese and Western traditions. In the 18th century, with the help of Western algebra, Mei Ch'ieh-ch'eng could understand again the ancient texts dealing with the method of the celestial unknown.

JAPAN IN THE 17TH CENTURY

The introduction of Chinese books. Very little is known of the history of Japanese mathematics before the 17th century. Through contacts with Chinese civilization, first mediated through Korea and then directly, there was, beginning in the seventh century, a flow of Chinese science to Japan. For example, the *Ten Classics of Mathematics* was introduced, along with the counting rods. Yet no Japanese book before the end of the 16th century is known to deal with mathematics. At that time another phase of importation began: the abacus and the *Systematic Treatise on Arithmetic* became known in Japan, where, however, they did not supplant the use of rods. Moreover, many books were taken from Korea to Japan, and perhaps in that way two Chinese books, the *Mathematical Treatise by Yang Hui* and the *Introduction to the Study of Mathematics* (written by Chu Shih-chieh in 1299), arrived in Japan. In those books, Japanese scholars could find the solution for systems of simultaneous linear equations, the search for the root of an equation according to the methods used in China in the 13th century, and applications of the method of the celestial unknown, although these were not easily understandable. Books about calendrical computations, containing mathematical knowledge, were also imported. Chinese mathematics greatly influenced the development of Japanese mathematics (for example, its algebraic orientation) and defined the context in which Japanese tradition opened to European mathematics.

Later, at the beginning of the Tokugawa period (1603–

1867), contacts with foreigners were limited to the trade with Chinese and Dutch boats through the harbour of Nagasaki. Some Chinese books, which may have contained Western knowledge, as well as Dutch books entered Japan secretly, but it is difficult to state how much, or what kind of, mathematical knowledge went through that channel.

The elaboration of Chinese methods. Although not the first mathematical book written in Japan, *Jingoki*, published in 1627 by Mitsuyoshi Yoshida, seems to be the first book that played an important part in the Japanese tradition. Inspired by the *Systematic Treatise on Arithmetic*, it described in Japanese the use of the soroban, an improvement of the Chinese abacus, and introduced some Chinese knowledge. Its many editions contributed to popularizing mathematics; most of the works on mathematics in Japan were written in Chinese and could not be widely read. In its enlarged edition of 1641, *Jingoki* introduced the method of performing computations with counting rods, which were no longer used in China. Moreover, Yoshida added "difficult problems" (*idai*), inspired by its Chinese source and left without solutions, which he recommended be posed to mathematicians. This initiated a tradition of challenges, reminiscent of those that took place in Europe during the Renaissance, that stimulated strongly the development of mathematics. In this context, in the 1650s, mathematicians, relying on counting-rod computations and looking for new methods of solution, began to understand by themselves the methods of Chinese algebra, as they were hinted at in the *Introduction to the Study of Mathematics*, which was reprinted in Japan in 1658.

These methods became a systematic tool for the solutions of any problem—what settled the algebraic framework for a mathematical development where equations and problems are essentially linked—and many books were published that introduced mathematicians to the method of finding the roots of equations or to the use of polynomial algebra to set down equations in order to solve problems. One of them, the *Kokon sanpoki* (1671), by Kazuyuki Sawaguchi, pointed out that "erroneous" problems could have more than one solution (equations could have more than one root) and left unanswered difficult problems involving simultaneous equations of the n th degree. The equations for their solutions were published in 1674 by Seki Kōwa (Takakazu), who was later referred to as the founder of the Japanese tradition of mathematics, called Wasan, that was based on Chinese mathematics. Seki founded a school of mathematics that became the most important one in Japan. As in the other schools, disciples had to keep the school methods secret, and only the best among them knew most of these methods. Only slowly did they publish their secrets, which hindered the free circulation of ideas (at this time, mathematics was widely practiced in Japan as a leisure activity) and which makes any attribution very difficult.

Explanations about how to use Seki's equations to derive Sawaguchi's problems were published in 1685 by one of his disciples, Takebe Katahiro. Seki had designed for this purpose a "literal" written algebra using characters, thus liberating mathematicians from counting rods; he kept for equations the positional notation with respect to one unknown, the coefficients being expressed in terms of other unknowns. In establishing equations among several unknowns for the solution of a problem, he had to introduce procedures equivalent to computations of determinants in order to eliminate unknowns between simultaneous equations. Further research elaborated these procedures.

Seki devised a classification of problems that amounted to a classification of equations, which took into consideration negative roots and multiple roots, noticed by Sawaguchi; for this purpose he adapted the Chinese algorithms from the 13th century. Seki and his disciples thus improved upon Chinese methods in many ways, opening new directions for the development of mathematics in Japan (for example, in their work on infinite series, the subject of research by contemporary European scientists as well).

(K.C.C.)

BIBLIOGRAPHY

General sources: Two standard texts are CARL B. BOYER, *A History of Mathematics*, 2nd ed. edited by UTA C. MERZBACH

Influence
of Chinese
texts

Wasan

(1989); and, on a more elementary level, HOWARD EVES, *An Introduction to the History of Mathematics*, 5th ed. (1983). Discussions of the mathematics of various periods may be found in O. NEUGEBAUER, *The Exact Sciences in Antiquity*, 2nd ed. (1957, reissued 1969); MORRIS KLINE, *Mathematical Thought from Ancient to Modern Times* (1972); and BARTEL L. VAN DER WAERDEN, *Science Awakening*, 4th ed. (1975; originally published in Dutch, 1950). See also KENNETH O. MAY, *Bibliography and Research Manual of the History of Mathematics* (1973); and JOSEPH W. DAUBEN, *The History of Mathematics from Antiquity to the Present: A Selective Bibliography* (1985). A good source for biographies of mathematicians is CHARLES COULSTON GILLISPIE (ed.), *Dictionary of Scientific Biography*, 16 vols. (1970–80, reissued 16 vol. in 8, 1981). Those wanting to study the writings of the mathematicians themselves will find the following source books useful: HENRIETTA O. MIDONICK (ed.), *The Treasury of Mathematics: A Collection of Source Material in Mathematics* (1965); JOHN FAUVEL and JEREMY GRAY (eds.), *The History of Mathematics: A Reader* (1987); D.J. STRUIK (ed.), *A Source Book in Mathematics, 1200–1800* (1969); and DAVID EUGENE SMITH, *A Source Book in Mathematics* (1929; reissued in 2 vol., 1959). A study of the development of numeric notation can be found in GEORGES IFRAH, *From One to Zero* (1985; originally published in French, 1981).

Mathematics in ancient Mesopotamia: Editions of mathematical tablets include O. NEUGEBAUER (ed. and trans.), *Mathematische Keilschrift-Texte*, 3 vol. (1935–37, reprinted 3 vol. in 2, 1973); and F. THUREAU-DANGIN (ed. and trans.), *Textes mathématiques babyloniens* (1938). O. NEUGEBAUER and A. SACHS, *Mathematical Cuneiform Texts* (1945), is the principal English edition of mathematical tablets. A brief look at Babylonian mathematics is contained in the first chapter of ASGER AABOE, *Episodes from the Early History of Mathematics* (1964), pp. 5–31.

Mathematics in ancient Egypt: Editions of the basic texts are T. ERIC PEET (ed. and trans.), *The Rhind Mathematical Papyrus* (1923, reprinted 1970); A.B. CHACE et al. (eds. and trans.), *The Rhind Mathematical Papyrus*, 2 vol. (1927–29); and W.W. STRUVE (V.V. STRUVE) (ed.), *Mathematischer papyrus des staatlichen Museums der schönen Künste in Moskau* (1930). A brief but useful summary appears in G.J. TOOMER, “Mathematics and Astronomy,” ch. 2 in J.R. HARRIS (ed.), *The Legacy of Egypt*, 2nd ed. (1971), pp. 27–54. For an extended account of Egyptian mathematics, see RICHARD J. GILLINGS, *Mathematics in the Time of the Pharaohs* (1972, reprinted 1982).

Greek mathematics: Critical editions of Greek mathematical texts include *The Thirteen Books of Euclid's Elements*, trans. by THOMAS L. HEATH, 2nd ed. rev., 3 vol. (1926, reprinted 1956); *The Works of Archimedes*, trans. by THOMAS L. HEATH (1897, reprinted 1953); E.J. DIJKSTERHUIS, *Archimedes*, trans. from Dutch (1956, reprinted 1987); THOMAS L. HEATH, *Apollonius of Perga: Treatise on Conic Sections* (1896, reissued 1961), and *Diophantus of Alexandria: A Study in the History of Greek Algebra*, 2nd ed. (1910, reprinted 1964); ROSHDI RASHED (trans.), *Les Arithmétiques* (1984–), of which vol. 3 and 4 contain Books IV–VII of Diophantus; and JACQUES SESIANO, *Books IV to VII of Diophantus' "Arithmetica" in the Arabic Translation Attributed to Qusṭā ibn Lūq* (1982). General surveys are THOMAS L. HEATH, *A History of Greek Mathematics*, 2 vol. (1921, reprinted 1981); JACOB KLEIN, *Greek Mathematical Thought and the Origin of Algebra* (1968; originally published in German, 1934); and WILBUR RICHARD KNORR, *The Ancient Tradition of Geometric Problems* (1986). Special topics are examined in O.A.W. DILKE, *Mathematics and Measurement* (1987); ÁRPÁD SZABÓ, *The Beginnings of Greek Mathematics* (1978; originally published in German, 1969); and WILBUR RICHARD KNORR, *The Evolution of the Euclidean Elements: A Study of the Theory of Incommensurable Magnitudes and Its Significance for Early Greek Geometry* (1975).

Mathematics in medieval Islām: Sources for Arabic mathematics include J.P. HOGENDIJK (ed. and trans.), *Ibn al-Haytham's Completion of the Conics*, trans. from Arabic (1985); MARTIN LEVEY and MARVIN PETRUCK (eds. and trans.), *Principles of Hindu Reckoning*, trans. from Arabic (1965), the only extant text of Kūshyār ibn Labbān's work; MARTIN LEVEY (ed. and trans.), *The Algebra of Abū Kāmil*, trans. from Arabic and Hebrew (1966), with a 13th-century Hebrew commentary by Mordecai Finzi; DAUD S. KASIR (ed. and trans.), *The Algebra of Omar Khayyam*, trans. from Arabic (1931, reprinted 1972); FREDERIC ROSEN (ed. and trans.), *The Algebra of Mohammed ben Musa*, trans. from Arabic (1831, reprinted 1986); and A.S. SAIDAN (ed. and trans.), *The Arithmetic of al-Uqlidisi*, trans. from Arabic (1978). Islāmic mathematics is examined in J.L. BERGGREN, *Episodes in the Mathematics of Medieval Islam* (1986); E.S. KENNEDY et al., *Studies in the Islamic Exact Sciences* (1983); and ROSHDI RASHIED, *Entre arithmétique et algèbre: recherches sur l'histoire des mathématiques arabes* (1984).

European mathematics during the Middle Ages and Renaissance: An overview is provided by MICHAEL S. MAHONEY, “Mathematics,” in DAVID C. LINDBERG (ed.), *Science in the Middle Ages* (1978), pp. 145–178. A.P. JUSCHKEWITSCH, *Geschichte der Mathematik im Mittelalter* (1964; originally published in Russian, 1961), pp. 326–434, is the definitive modern work. Other sources include ALEXANDER MURRAY, *Reason and Society in the Middle Ages* (1978, reprinted 1985), ch. 6–8; GEORGE SARTON, *Introduction to the History of Science*, vol. 2, *From Rabbi Ben Ezra to Roger Bacon*, 2 parts (1931, reprinted 1975), and vol. 3, *Science and Learning in the Fourteenth Century*, 2 parts (1947–48, reprinted 1975); and, on a more advanced level, EDWARD GRANT and JOHN E. MURDOCH (eds.), *Mathematics and Its Applications to Science and Natural Philosophy in the Middle Ages* (1987). For the Renaissance, see PAUL LAWRENCE ROSE, *The Italian Renaissance of Mathematics: Studies on Humanists and Mathematicians from Petrarch to Galileo* (1975).

Mathematics in the 17th and 18th centuries: An overview of this period is contained in DEREK THOMAS WHITESIDE, “Patterns of Mathematical Thought in the Later Seventeenth Century.” *Archive for History of Exact Sciences*, 1(3):179–388 (1961). Specific topics are examined in MARGARET E. BARON, *The Origins of the Infinitesimal Calculus* (1969, reprinted 1987); ROBERTO BONOLA, *Non-Euclidean Geometry: A Critical and Historical Study of Its Development* (1955; originally published in Italian, 1912); CARL B. BOYER, *The Concepts of the Calculus: A Critical and Historical Discussion of the Derivative and the Integral* (1930, reissued with the title *The History of the Calculus and Its Conceptual Development*, 1949, reprinted 1959); HERMAN GOLDSTINE, *A History of Numerical Analysis from the 16th Through the 19th Century* (1977); JUDITH V. GRABINER, *The Origins of Cauchy's Rigorous Calculus* (1981); I. GRATTAN-GUINNESS, *The Development of the Foundations of Mathematical Analysis from Euler to Riemann* (1970); ROGER HAHN, *The Anatomy of a Scientific Institution: The Paris Academy of Sciences, 1666–1803* (1971); and LUBOŠ NOVÝ, *Origins of Modern Algebra*, trans. from Czech (1973).

Mathematics in the 19th and 20th centuries: Surveys include HERBERT MEHRTENS, HENK BOS, and IVO SCHNEIDER (eds.), *Social History of Nineteenth Century Mathematics* (1981); WILLIAM ASPRAY and PHILIP KITCHER (eds.), *History and Philosophy of Modern Mathematics* (1988); and KEITH DEVLIN, *Mathematics: The New Golden Age* (1988). Special topics are examined in UMBERTO BOTTAZZINI, *The Higher Calculus: A History of Real and Complex Analysis from Euler to Weierstrass* (1986; originally published in Italian, 1981); JULIAN LOWELL COOLIDGE, *A History of Geometrical Methods* (1940, reissued 1963); JOSEPH WARREN DAUBEN, *Georg Cantor: His Mathematics and Philosophy of the Infinite* (1979); HAROLD M. EDWARDS, *Fermat's Last Theorem: A Genetic Introduction to Algebraic Number Theory* (1977); I. GRATTAN-GUINNESS (ed.), *From the Calculus to Set Theory, 1630–1910: An Introductory History* (1980); JEREMY GRAY, *Ideas of Space: Euclidian, Non-Euclidean, and Relativistic* (1979); THOMAS HAWKINS, *Lebesgue's Theory of Integration: Its Origins and Development*, 3rd ed. (1979); JESPER LÜTZEN, *The Prehistory of the Theory of Distributions* (1982); and MICHAEL MONASTYRSKY, *Riemann, Topology, and Physics*, trans. from Russian (1987).

Mathematics in China and Japan: Chinese mathematics is discussed in JOSEPH NEEDHAM, *Science and Civilisation in China*, vol. 3, *Mathematics and the Sciences of the Heavens and the Earth* (1959, reprinted 1970), pp. 1–168; *Ancient China's Technology and Science* (1983), a group of papers prepared by the Institute of the History of Natural Sciences, Chinese Academy of Science, in Peking; YAN LI (YEN LI) and SHIRAN DU (SHIH-JAN TU), *Chinese Mathematics: A Concise History*, trans. from Chinese (1987); ULRICH LIBBRECHT, *Chinese Mathematics in the Thirteenth Century* (1973); and LAY YONG LAM, *A Critical Study of the Yang Hui Suan Fa: A Thirteenth-Century Chinese Mathematical Treatise*, trans. from Chinese (1977). Useful journal articles include DONALD BLACKMORE WAGNER, “An Early Chinese Derivation of the Volume of a Pyramid: Liu Hui, Third Century A.D.,” *Historia Mathematica*, 6(2):164–188 (May 1979), and “Liu Hui and Tsu Keng-Chih on the Volume of a Sphere,” *Chinese Science*, 3:59–79 (1978). Besides *Historia Mathematica* (quarterly) and *Chinese Science* (irregular), many papers on Chinese mathematics may be found in *Archive for History of Exact Sciences* (8/yr.).

Overviews of Japanese mathematics include DAVID EUGENE SMITH and YOSHIO MIKAMI, *A History of Japanese Mathematics* (1914); YOSHIO MIKAMI, *The Development of Mathematics in China and Japan*, 2nd ed. (1974); and SHIGERU NAKAYAMA, “Japanese Scientific Thought,” in CHARLES COULSTON GILLISPIE (ed.), *Dictionary of Scientific Biography*, vol. 15 (1978), pp. 728–758. *Historia Scientiarum: International Journal of the History of Science Society of Japan* (annual), contains many papers on Japanese mathematics.

(J.L.B./W.R.K./Mc.F./C.G.F./J.J.G./K.C.C.)

Matter: Its Properties, States, Varieties, and Behaviour

The tangible universe—that is, everything that has mass and occupies space—is made of matter. Because it is difficult to identify anything that is not matter, it is more meaningful to consider the specific characteristics of matter than to attempt to provide a rigorous definition.

This article deals with those important forms of matter typically encountered—such as solids, liquids, and gases—for which the basic building blocks are atoms. (Other types of matter do exist, such as the ultradense nuclear matter of neutron stars and the isolated subatomic particles in interstellar space.) There are more than 100 distinctly different types of atoms, corresponding to the different chemical elements. Each atom consists of negatively charged point-like electrons swarming around a small positively charged nucleus made up of protons and neutrons. The number of electrons ranges from 1 (for hydrogen) to 92 (for uranium, the heaviest known naturally occurring element) to 110 (for the heaviest element thus far known). The diameter of a typical atom, corresponding to the size of the electron swarm, is about 100 million times smaller than one inch. The nucleus within the atom is considerably smaller; it has a diameter that is roughly 100,000 times smaller than the atomic diameter. The electron is believed to be an undividable elementary building block of matter. The protons and neutrons that make up the nucleus are themselves made up of more elementary building blocks (called quarks). This internal structure of the nucleus does not enter into the properties of solids, liquids, and gases, however. It is the interactions between the electron swarms that are responsible for binding atoms together into small

chemically bonded groups of atoms called molecules and into the large groups of bonded atoms that constitute solids. In solids the atoms are in intimate contact, and chemical bonds keep them locked into a rigid structure; in liquids they are close together but have the freedom to maneuver around each other; and in gases they are far apart (relative to their size).

This article describes the properties of the different phases of matter: solid, liquid, and gas. A given material can exist, depending on the conditions, in different phases, and so changes from one phase to another are discussed. Within the solid state, there are important distinctions (especially in the way that the atoms are arranged) between crystalline solids, quasicrystals, and amorphous solids. Liquid crystals are complex anisotropic liquids that exhibit unique properties, as do plasmas, which are electrically conducting collections of ionized particles. In the plasma state (found in the Sun, for example) high temperatures cause electrons and nuclei to separate and exhibit pronounced electrical behaviour. Another form of matter consists of atomic or molecular clusters that are too small to act as bulk solids and yet behave quite differently from individual atoms or molecules. The interesting properties of materials at low temperature and at high pressure are described as well. For additional information about the basic constituents and structure of matter, see the articles SUBATOMIC PARTICLES, ATOMS, and CHEMICAL ELEMENTS. (R.Z.)

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, sections 124 and 125, and the *Index*.

This article is divided into the following sections:

-
- Phase changes 613
 - General considerations
 - System variables
 - Applications to petrology
 - Solid state 615
 - Crystalline solids 615
 - Classification
 - Structure
 - Types of bonds
 - Crystal growth
 - Electric properties
 - Magnetism
 - Quasicrystals 628
 - Structure and symmetry
 - Properties
 - Liquid crystals 632
 - Structure and symmetry
 - Optical properties
 - Amorphous solids 636
 - Distinction between crystalline and amorphous solids
 - Preparation of amorphous solids
 - Atomic-scale structure
 - Properties of oxide glasses
 - Properties and applications of amorphous solids
 - Liquid state 643
 - Behaviour of pure liquids 644
 - Molecular structure of liquids
 - Speed of sound and electric properties
 - Solutions and solubilities 646
 - Classes of solutions
 - Properties of solutions
 - Thermodynamics and intermolecular forces in solutions
 - Theories of solutions
 - Solubilities of solids and gases
 - Gaseous state 656
 - Structure 656
 - Kinetic-molecular picture
 - Numerical magnitudes
 - Free-molecule gas
 - Continuity of gaseous and liquid states
 - Behaviour and properties 659
 - Equilibrium properties
 - Transport properties
 - Kinetic theory of gases 661
 - Ideal gas
 - Deviations from the ideal model
 - Plasma state 665
 - Basic plasma physics 666
 - Plasma formation
 - Methods of describing plasma phenomena
 - Determination of plasma variables
 - Waves in plasmas
 - Applications of plasmas
 - Natural plasmas 669
 - Extraterrestrial forms
 - Solar-terrestrial forms
 - Clusters 672
 - Comparison with other forms of matter 673
 - Methods of study 674
 - Structure and properties 675
 - Structure
 - Physical properties
 - Chemical properties
 - Low-temperature phenomena 678
 - Superconductivity 678
 - Thermal properties of superconductors
 - Magnetic and electromagnetic properties of superconductors
 - Higher-temperature superconductivity
 - Superfluidity 682
 - High-pressure phenomena 683
 - Producing high pressure 684
 - Physical and chemical effects of high pressure 685
 - Applications 687
 - Bibliography 688
-

PHASE CHANGES

GENERAL CONSIDERATIONS

A system is a portion of the universe that has been chosen for studying the changes that take place within it in response to varying conditions. A system may be complex, such as a planet, or relatively simple, as the liquid within a glass. Those portions of a system that are physically distinct and mechanically separable from other portions of the system are called phases. When a phase in one form is altered to another form, a phase change is said to have occurred.

Phases within a system exist in a gaseous, liquid, or solid state. Solids are characterized by strong atomic bonding and high viscosity, resulting in a rigid shape. Most solids are crystalline, inasmuch as they have a three-dimensional periodic atomic arrangement; some solids (such as glass) lack this periodic arrangement and are noncrystalline, or amorphous. Gases consist of weakly bonded atoms with no long-range periodicity; gases expand to fill any available space. Liquids have properties intermediate between those of solids and gases. The molecules of a liquid are condensed like those of a solid. Liquids have a definite volume, but their low viscosity enables them to change shape as a function of time. The matter within a system may consist of more than one solid or liquid phase, but a system can contain only a single gas phase, which must be of homogeneous composition because the molecules of gases mix completely in all proportions.

SYSTEM VARIABLES

Systems respond to changes in pressure, temperature, and chemical composition, and, as this happens, phases may be created, eliminated, or altered in composition. For example, an increase in pressure may cause a low-density liquid to convert to a denser solid, while an increase in temperature may cause a solid to melt. A change of composition might result in the compositional modification of a preexisting phase or in the gain or loss of a phase.

The classification and limitations of phase changes are described by the phase rule, as proposed by the American chemist J. Willard Gibbs in 1876 and based on a rigorous thermodynamic relationship. The phase rule is commonly given in the form $P + F = C + 2$. The term P refers to the number of phases that are present within the system, and C is the minimum number of independent chemical components that are necessary to describe the composition of all phases within the system. The term F , called the variance, or degrees of freedom, describes the minimum number of variables that must be fixed in order to define a particular condition of the system.

Phase relations are commonly described graphically in terms of phase diagrams (see Figure 1). Each point within the diagram indicates a particular combination of pressure and temperature, as well as the phase or phases that exist stably at this pressure and temperature. All phases in Figure 1 have the same composition—that of silicon dioxide, SiO_2 . The diagram is a representation of a one-component (unary) system, in contrast to a two-component (binary), three-component (ternary), or four-component (quaternary) system. The phases coesite, low quartz, high quartz, tridymite, and cristobalite are solid phases composed of silicon dioxide; each has its own atomic arrangement and distinctive set of physical and chemical properties. The most common form of quartz (found in beach sands and granites) is low quartz. The region labeled anhydrous melt consists of silicon dioxide liquid.

Different portions of the silicon dioxide system may be examined in terms of the phase rule. At point A a single solid phase exists—low quartz. Substituting the appropriate values into the phase rule $P + F = C + 2$ yields $1 + F = 1 + 2$, so $F = 2$. For point A (or any point in which only a single phase is stable) the system is divariant—*i.e.*, two degrees of freedom exist. Thus, the two variables (pressure and temperature) can be changed independently, and the same phase assemblage continues to exist.

Point B is located on the boundary curve between the stability fields of low quartz and high quartz. At all points along this curve, these two phases coexist. Substituting values in the phase rule ($2 + F = 1 + 2$) will cause a variance of 1 to be obtained. This indicates that one independent variable can be changed such that the same pair of phases will be retained. A second variable must be changed to conform to the first in order for the phase assemblage to remain on the boundary between low and high quartz. The same result holds for the other boundary curves in this system.

Point C is located at a triple point, a condition in which three stability fields intersect. The phase rule ($3 + F = 1 + 2$) indicates that the variance is 0. Point C is therefore an invariant point; a change in either pressure or temperature results in the loss of one or more phases. The phase rule also reveals that no more than three phases can stably coexist in a one-component system because additional phases would lead to negative variance.

From W. G. Ernst, *Petrologic Phase Equilibria*, copyright © 1976 W. H. Freeman and Company, used by permission, after F. R. Boyd and J. L. England, "The Quartz-Coesite Transition," *Journal of Geophysical Research*, vol. 65, no. 2, February 1960

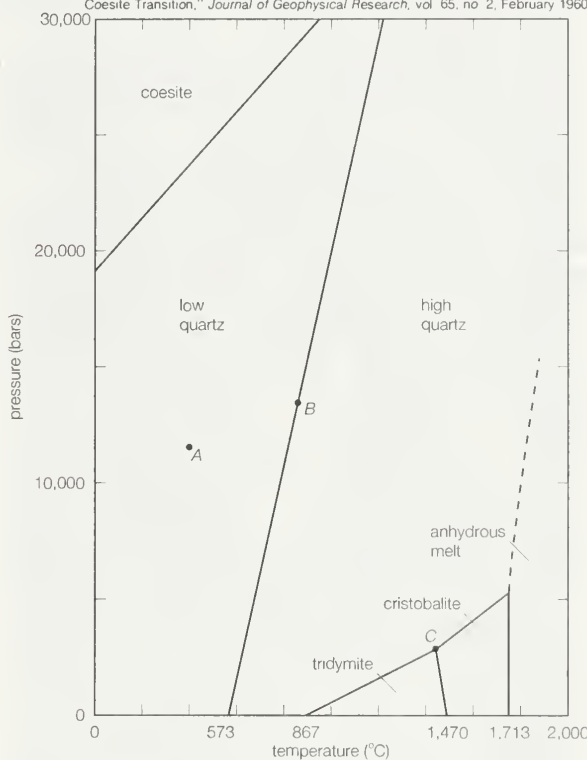


Figure 1: Phase diagram for the unary system SiO_2 .

Consider the binary system (Figure 2) that describes the freezing and melting of the minerals sphene (CaSiTiO_5) or titanite, and anorthite feldspar ($\text{CaAl}_2\text{Si}_2\text{O}_8$). The melt can range in composition from pure CaSiTiO_5 to pure $\text{CaAl}_2\text{Si}_2\text{O}_8$, but the solids show no compositional substitution. All phases therefore have the composition of CaSiTiO_5 , $\text{CaAl}_2\text{Si}_2\text{O}_8$, or a liquid mixture of the two. The system in the figure has been examined at atmospheric pressure; because the pressure variable is fixed, the phase rule is expressed as $P + F = C + 1$. In this form it is called the condensed phase rule, for any gas phase is either condensed to a liquid or is present in negligible amounts. The phase diagram shows a vertical temperature coordinate and a horizontal compositional coordinate (ranging from pure CaSiTiO_5 at the left to pure $\text{CaAl}_2\text{Si}_2\text{O}_8$ at the right).

The phase fields (separated by the solid curves) contain either one or two phases. Any point in a one-phase field corresponds to a single phase whose composition is indicated directly below on the horizontal axis. For example, point A presents a liquid whose composition is 70 percent $\text{CaAl}_2\text{Si}_2\text{O}_8$ and 30 percent CaSiTiO_5 . The compositions of phases in a two-phase field are determined by construction

The phase rule

The unary system SiO_2

The binary sphene-anorthite system

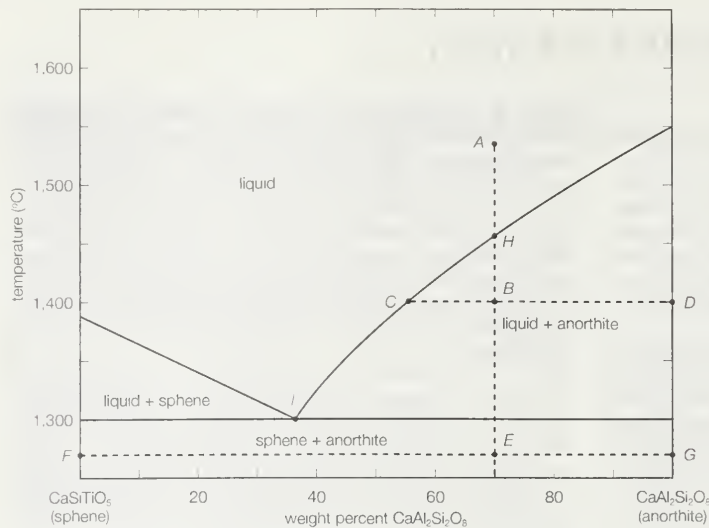


Figure 2: Phase diagram for the binary system CaSiTiO_5 - $\text{CaAl}_2\text{Si}_2\text{O}_8$.

From E. G. Ehlers, *The Interpretation of Geological Phase Diagrams*, © 1987 Dover Publications, Inc., after A. T. Prince, "The System Albite-Anorthite-Sphene," *The Journal of Geology*, vol. 51, no. 1, 1943, © The University of Chicago Press, publisher

of a horizontal (constant-temperature) line from the point of interest to the extremities of the two-phase field. Thus, a sample with composition *B* consists of liquid *C* (43 percent CaSiTiO_5 and 57 percent $\text{CaAl}_2\text{Si}_2\text{O}_8$) and solid anorthite *D*. A sample at point *E* at a lower temperature consists of the solids sphene (*F*) and anorthite (*G*).

Liquid $\text{CaAl}_2\text{Si}_2\text{O}_8$ cools to produce solid anorthite at 1,550° C, whereas liquid CaSiTiO_5 cools to produce solid sphene at 1,390° C. If the batch were a mixture of the two components, the freezing temperature of each of these minerals would be depressed. In a melt consisting of a single component, such as CaSiTiO_5 , all atoms could add to sphene nuclei to form crystals of sphene. If, however, the melt contained 30 percent $\text{CaAl}_2\text{Si}_2\text{O}_8$, the rate of formation of sphene nuclei would be decreased, as 30 percent of the melt could not contribute to their formation. In order to increase the rate of formation of sphene nuclei and promote crystallization, the temperature of the melt must be decreased below the freezing point of pure CaSiTiO_5 . When cooled, liquid *A* does not begin crystallization until temperature *II* is reached. Pure anorthite crystals precipitate from the melt. Depletion of $\text{CaAl}_2\text{Si}_2\text{O}_8$ from the melt causes the melt composition to become relatively enriched in CaSiTiO_5 , with consequent additional depression of the anorthite freezing point. As freezing continues, the liquid composition changes until the minimum point is reached at *I*. This point is called the eutectic. It is the lowest temperature at which a liquid can exist in this system. At the eutectic, both anorthite and sphene crystallize together at a fixed temperature and in a fixed ratio until the remaining liquid is consumed. All intermediate liquid compositions migrate during crystallization

to the eutectic. The melting sequence of sphene-anorthite mixtures is exactly the opposite of the freezing sequence (*i.e.*, melting of any anorthite-sphene mixture begins at the eutectic).

Depression of the freezing point of a compound by the addition of a second component is common in both binary and more complex systems. This usually occurs when the solid phases either have a fixed composition or show limited solid solution. Common examples are the mixing of ice and salt (NaCl) or the use of ethylene glycol (anti-freeze) to depress the freezing point of water (see Figure 3).

APPLICATIONS TO PETROLOGY

Systematic investigation of the phase changes of the more common anhydrous mineral groups was initiated by the Canadian-born American petrologist Norman L. Bowen and his coworkers at the Geophysical Laboratory of the Carnegie Institution of Washington, D.C., in the early 20th century. This work was generally limited to systems at atmospheric pressure. Subsequent advances in technology have permitted the examination of rock systems in the presence of water pressure and ultrahigh confining pressures. Materials can now be systematically examined under conditions that range from those at the Earth's surface to those simulating conditions that exist at the core. This has led to a vast increase in knowledge about the conditions of formation of both igneous and metamorphic rocks. Synthetic equivalents of almost every mineral or rock system can now be produced in the laboratory. Even gemstones such as diamonds are routinely synthesized.

Typical of the data now available are the freezing-melting curves (Figure 4) of the common volcanic rock basalt (and its coarse-grained equivalent, gabbro). Figure 4A shows the crystallization range (shaded) for basaltic melts as a function of lithostatic pressure; this pressure is due to depth of burial. The two short lines show the approximate position of a transition region between gabbro and its denser solid equivalent, eclogite (a sodium-pyroxene + garnet rock). The melting curves have a positive slope, as the solids are denser than their equivalent melts and are thus favoured (enlarged) with increasing pressure.

In the presence of water pressure ($P_{\text{H}_2\text{O}}$), the freezing-melting curves are depressed (Figure 4B) because the water acts as another component. The slope of the curves is also influenced by the presence of a hydrous solid phase, hornblende, whose approximate stability field is indicated by the dashed line. The changes in liquid composition and crystallization sequences have been determined. Similar information is available for most common igneous rocks.

In 1915 the Finnish petrologist Pentti E. Eskola set up a classification scheme for metamorphic rocks that was based on metamorphic facies. Each facies was defined by the presence of one or more common mineral assemblages.

Freezing-point depression

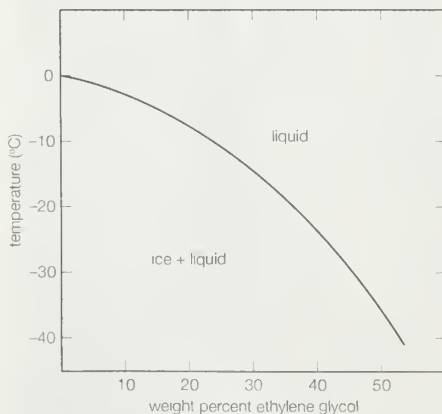


Figure 3: The depression of the freezing point of water as a function of the addition of ethylene glycol.

Rock formation

The stability limits of these assemblages subsequently have been determined by laboratory studies. As a result, placing a metamorphic rock within a particular facies indicates the broad pressure-temperature region in which the rock formed. (See MINERALS AND ROCKS: *Metamorphic rocks: Metamorphic facies* for the pressure-temperature regions of the major metamorphic facies.) For example, a rock containing sodium-rich pyroxene and garnet is placed within the eclogite facies, which indicates that it formed at

pressures greater than about 12 kilobars and temperatures above approximately 600° C. Rocks in the blueschist facies contain the blue amphibole glaucophane; such rocks are stable at high pressures and relatively low temperatures.

A large variety of schemes are available to provide more detailed information on the temperatures and pressures of formation of both igneous and metamorphic rocks. These may use phase relations, stable isotopes, or the compositions of coexisting mineral pairs. (E.G.E.)

SOLID STATE

Crystalline solids

CLASSIFICATION

The definition of a solid appears obvious; a solid is generally thought of as being hard and firm. Upon inspection, however, the definition becomes less straightforward. A cube of butter, for example, is hard after being stored in a refrigerator and is clearly a solid. After remaining on the kitchen counter for a day, the same cube becomes quite soft, and it is unclear if the butter should still be considered a solid. Many crystals behave like butter in that they are hard at low temperatures but soft at higher temperatures. They are called solids at all temperatures below their melting point. A possible definition of a solid is an object that retains its shape if left undisturbed. The pertinent issue is how long the object keeps its shape. A highly viscous fluid retains its shape for an hour but not a year. A solid must keep its shape longer than that.

The basic units of solids are either atoms or atoms that have combined into molecules. The electrons of an atom move in orbits that form a shell structure around the

nucleus. The shells are filled in a systematic order, with each shell accommodating only a small number of electrons. Different atoms have different numbers of electrons, which are distributed in a characteristic electronic structure of filled and partially filled shells. The arrangement of an atom's electrons determines its chemical properties. The properties of solids are usually predictable from the properties of their constituent atoms and molecules, and the different shell structures of atoms are therefore responsible for the diversity of solids.

All occupied shells of the argon (Ar) atom, for example, are filled, resulting in a spherical atomic shape. In solid argon the atoms are arranged according to the closest packing of these spheres. The iron (Fe) atom, in contrast, has one electron shell that is only partially filled, giving the atom a net magnetic moment. Thus, crystalline iron is a magnet. The covalent bond between two carbon (C) atoms is the strongest bond found in nature. This strong bond is responsible for making diamond the hardest solid.

A solid is crystalline if it has long-range order. Once the positions of an atom and its neighbours are known at one point, the place of each atom is known precisely throughout the crystal. Most liquids lack long-range order, although many have short-range order. Short range is defined as the first- or second-nearest neighbours of an atom. In many liquids the first-neighbour atoms are arranged in the same structure as in the corresponding solid phase. At distances that are many atoms away, however, the positions of the atoms become uncorrelated. These fluids, such as water, have short-range order but lack long-range order. Certain liquids may have short-range order in one direction and long-range order in another direction; these special substances are called liquid crystals. Solid crystals have both short-range order and long-range order.

Solids that have short-range order but lack long-range order are called amorphous. Almost any material can be made amorphous by rapid solidification from the melt (molten state). This condition is unstable, and the solid will crystallize in time. If the timescale for crystallization is years, then the amorphous state appears stable. Glasses are an example of amorphous solids; they are discussed below in the section *Amorphous solids*. In crystalline silicon (Si) each atom is tetrahedrally bonded to four neighbours. In amorphous silicon (a-Si) the same short-range order exists, but the bond directions become changed at distances farther away from any atom. Amorphous silicon is a type of glass. Quasicrystals are another type of solid that lack long-range order (see below *Quasicrystals*).

Most solid materials found in nature exist in polycrystalline form rather than as a single crystal. They are actually composed of millions of grains (small crystals) packed together to fill all space. Each individual grain has a different orientation than its neighbours. Although long-range order exists within one grain, at the boundary between grains, the ordering changes direction. A typical piece of iron or copper (Cu) is polycrystalline. Single crystals of metals are soft and malleable, while polycrystalline metals are harder and stronger and are more useful industrially. Most polycrystalline materials can be made into large single crystals after extended heat treatment. In the past blacksmiths would heat a piece of metal to make it malleable; heat makes a few grains grow large by incorporating smaller ones. The smiths would bend the softened metal into shape and then pound it awhile; the pounding would make it polycrystalline again, increasing its strength.

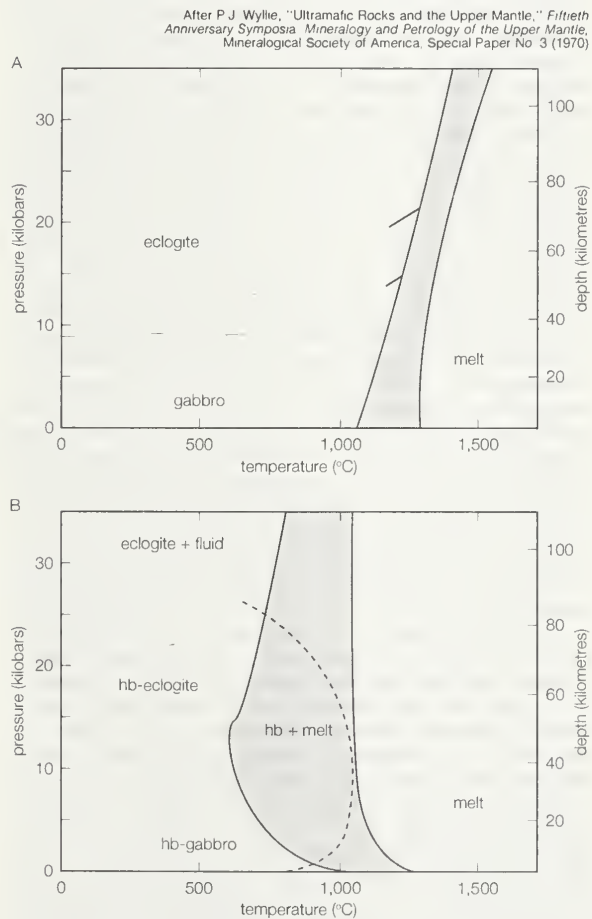


Figure 4: Freezing-melting curves. (A) The crystallization range of basaltic (gabbroic) bulk composition (shaded) as a function of lithostatic (burial) pressure. The two short lines indicate the approximate position of a transition between gabbro and its denser equivalent, eclogite. (B) The crystallization range (shaded) of gabbroic bulk composition as a function of water pressure. The abbreviation "hb" denotes hornblende.

Long-range order

Categories of crystals

Crystals are classified in general categories, such as insulators, metals, semiconductors, and molecular solids. A single crystal of an insulator is usually transparent and resembles a piece of glass. Metals are shiny unless they have rusted. Semiconductors are sometimes shiny and sometimes transparent but are never rusty. Many crystals can be classified as a single type of solid, while others have intermediate behaviour. Cadmium sulfide (CdS) can be prepared in pure form and is an excellent insulator; when impurities are added to cadmium sulfide, it becomes an interesting semiconductor. Bismuth (Bi) appears to be a metal, but the number of electrons available for electrical conduction is similar to that of semiconductors. In fact, bismuth is called a semimetal. Molecular solids are usually crystals formed from molecules or polymers. They can be insulating, semiconducting, or metallic, depending on the type of molecules in the crystal. New molecules are continuously being synthesized, and many are made into crystals. The number of different crystals is enormous.

STRUCTURE

Crystals can be grown under moderate conditions from all 92 naturally occurring elements except helium, and helium can be crystallized at low temperatures by using 25 atmospheres of pressure. Binary crystals are composed of two elements. There are thousands of binary crystals; some examples are sodium chloride (NaCl), alumina (Al_2O_3), and ice (H_2O). Crystals can also be formed with three or more elements.

A basic concept in crystal structures is the unit cell. It is the smallest unit of volume that permits identical cells to be stacked together to fill all space. By repeating the pattern of the unit cell over and over in all directions, the entire crystal lattice can be constructed. A cube is the simplest example of a unit cell. Two other examples are shown in Figure 5. The first is the unit cell for a face-centred cubic lattice, and the second is for a body-centred cubic lattice. These structures are explained in the following paragraphs. There are only a few different unit-cell shapes, so many different crystals share a single unit-cell type. An important characteristic of a unit cell is the number of atoms it contains. The total number of atoms in the entire crystal is the number in each cell multiplied by the number of unit cells. Copper and aluminum (Al)

each have one atom per unit cell, while zinc (Zn) and sodium chloride have two. Most crystals have only a few atoms per unit cell, but there are some exceptions. Crystals of polymers, for example, have thousands of atoms in each unit cell.

The elements are found in a variety of crystal packing arrangements. The most common lattice structures for metals are those obtained by stacking the atomic spheres into the most compact arrangement. There are two such possible periodic arrangements. In each, the first layer has the atoms packed into a plane-triangular lattice in which every atom has six immediate neighbours. Figure 6 shows this arrangement for the atoms labeled *A*. The second layer is shaded in the figure. It has the same plane-triangular structure; the atoms sit in the holes formed by the first layer. The first layer has two equivalent sets of holes, but the atoms of the second layer can occupy only one set. The third layer, labeled *C*, has the same structure, but there are two choices for selecting the holes that the atoms will occupy. The third layer can be placed over the atoms of the first layer, generating an alternate layer sequence *ABABAB* . . . , which is called the hexagonal-closest-packed (hcp) structure. Cadmium and zinc crystallize with this structure. The second possibility is to place the atoms of the third layer over those of neither of the first two but instead over the set of holes in the first layer that remains unoccupied. The fourth layer is placed over the first, and so there is a three-layer repetition *ABCABCABC* . . . , which is called the face-centred cubic (fcc), or cubic-closest-packed, lattice. Copper, silver (Ag), and gold (Au) crystallize in fcc lattices. In the hcp and the fcc structures the spheres fill 74 percent of the volume, which represents the closest possible packing of spheres. Each atom has 12 neighbours. The number of atoms in a unit cell is two for hcp structures and one for fcc. There are 32 metals that have the hcp lattice and 26 with the fcc. Another possible arrangement is the body-centred cubic (bcc) lattice, in which each atom has eight neighbours arranged at the corners of a cube. Figure 7A shows the cesium chloride (CsCl) structure, which is a cubic arrangement. If all atoms in this structure are of the same species, it is a bcc lattice. The spheres occupy 68 percent of the volume. There are 23 metals with the bcc arrangement. The sum of these three numbers (32 + 26 + 23) exceeds the number of elements that form metals (63), since some elements are found in two or three of these structures.

The fcc structure is also found for crystals of the rare gas solids neon (Ne), argon (Ar), krypton (Kr), and xenon (Xe). Their melting temperatures at atmospheric pressure are: Ne, 24.6 K; Ar, 83.8 K; Kr, 115.8 K; and Xe, 161.4 K.

The elements in the fourth row of the periodic table—carbon, silicon, germanium (Ge), and α -tin (α -Sn)—prefer covalent bonding. Carbon has several possible crystal structures. Each atom in the covalent bond has four first-neighbours, which are at the corners of a tetrahedron. This arrangement is called the diamond lattice and is shown in Figure 7C. There are two atoms in a unit cell, which is fcc. Large crystals of diamond are valuable gemstones. The crystal has other interesting properties; it has the highest sound velocity of any solid and is the best conductor of heat. Besides diamond, the other common form of carbon is graphite, which is a layered material. Each carbon atom has three coplanar near neighbours, forming an arrangement called the honeycomb lattice. Three-dimensional graphite crystals are obtained by stacking similar layers.

Another form of crystalline carbon is based on a molecule with 60 carbon atoms called buckminsterfullerene (C_{60}). The molecular shape is spherical. Each carbon is bonded to three neighbours, as in graphite, and the spherical shape is achieved by a mixture of 12 rings with five sides and 20 rings with six sides. Similar structures were first visualized by the American architect R. Buckminster Fuller for geodesic domes. The C_{60} molecules, also called buckyballs, are quite strong and almost incompressible. Crystals are formed such that the balls are arranged in an fcc lattice with a one-nanometre spacing between the centres of adjacent balls. The similar C_{70} molecule has the shape of a rugby ball; C_{70} molecules also form an fcc crystal when stacked together. The solid fullerenes form

Arrangements of atoms in crystals

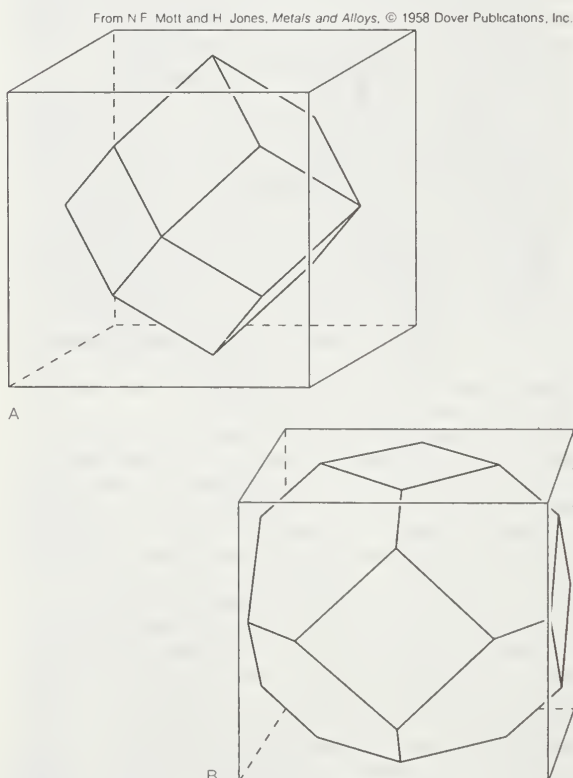


Figure 5: Unit cells for (A) face-centred and (B) body-centred cubic lattices.

Structures of carbon

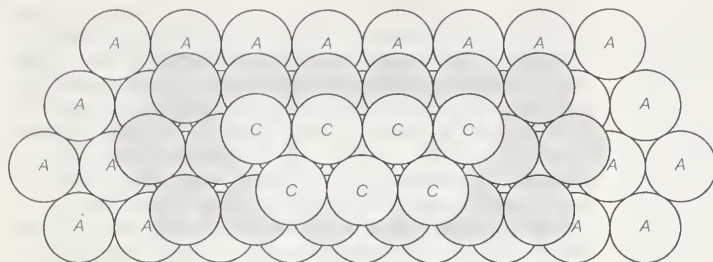


Figure 6: Stacking of spheres in closest-packed arrangements. The atoms of the second layer (shaded) sit in the holes formed by the atoms of the first layer (labeled A). If the atoms of the third layer are placed in positions labeled C and the structure is repeated (ABCABC...), the face-centred cubic lattice is formed. If the third layer is placed directly over the atoms of the first layer (A) and the structure is repeated (ABABAB...), the hexagonal closest-packed arrangement is obtained.

molecular crystals, with weak binding—provided by van der Waals interactions—between the molecules.

Many elements form diatomic gases: hydrogen (H), oxygen (O), nitrogen (N), fluorine (F), chlorine (Cl), bromine (Br), and iodine (I). When cooled to low temperature, they form solids of diatomic molecules. Nitrogen has the hcp structure, while oxygen has a more complex structure.

The most interesting crystal structures are those of elements that are neither metallic, covalent, nor diatomic. Although boron (B) and sulfur (S) have several different crystal structures, each has one arrangement in which it is usually found. Twelve boron atoms form a molecule in the shape of an icosahedron (Figure 8). Crystals are formed by stacking the molecules. The β -rhombohedral structure of boron has seven of these icosahedral molecules in each unit cell, giving a total of 84 atoms. Molecules of sulfur are usually arranged in rings; the most common ring has eight atoms. The typical structure is α -sulfur, which has 16 molecules per unit cell, or 128 atoms. In the common crystals of selenium (Se) and tellurium (Te), the atoms are arranged in helical chains, which stack like cordwood. However, selenium also makes eight-atom rings, similar to sulfur, and forms crystals from them. Sulfur also makes helical chains, similar to selenium, and stacks them together into crystals.

Binary crystals are found in many structures. Some pairs of elements form more than one structure. At room temperature, cadmium sulfide may crystallize either in the zinc blende or wurtzite structure. Alumina also has two possible structures at room temperature, α -alumina (corundum) and β -alumina. Other binary crystals exhibit different structures at different temperatures. Among the most complex crystals are those of silicon dioxide (SiO_2), which has seven different structures at various temperatures and pressures; the most common of these structures is quartz. Some pairs of elements form several different crystals in which the ions have different chemical valences. Cadmium (Cd) and phosphorus (P) form the crystals Cd_3P_2 , CdP_2 , CdP_4 , Cd_7P_{10} , and Cd_8P_7 . Only in the first case are the ions assigned the expected chemical valences of Cd^{2+} and P^{3-} .

Among the binary crystals, the easiest structures to visualize are those with equal numbers of the two types of atoms. The structure of sodium chloride is based on a cube. To construct the lattice, the sodium and chlorine atoms are placed on alternate corners of a cube, and the structure is repeated (Figure 7B). The structure of the sodium atoms alone, or the chlorine atoms alone, is fcc and defines the unit cell. The sodium chloride structure thus is made up of two interpenetrating fcc lattices. The cesium chloride lattice (Figure 7A) is based on the bcc structure; every other atom is cesium or chlorine. In this case, the unit cell is a cube. The third important structure for AB (binary) lattices is zinc blende (Figure 7D). It is based on the diamond structure, where every other atom is A or B. Many binary semiconductors have this structure, including those with one atom from the third (boron, aluminum, gallium [Ga], or indium [In]) and one from the fifth (nitrogen, phosphorus, arsenic [As], or anti-

mony [Sb]) column of the periodic table (GaAs, InP, etc.). Most of the chalcogenides (O, S, Se, Te) of cadmium and zinc (CdTe, ZnSe, ZnTe, etc.) also have the zinc blende structure. The mineral zinc blende is ZnS; its unit cell is also fcc. The wurtzite structure is based on the hcp lattice, where every other atom is A or B. These four structures comprise most of the binary crystals with equal numbers of cations and anions.

The fullerene molecule forms binary crystals $M_x\text{C}_{60}$ with alkali atoms, where M is potassium (K), rubidium (Rb), or cesium (Cs). The fullerene molecules retain their spherical shape, and the alkali atoms sit between them. The subscript x can take on several values. A compound with $x=6$ (e.g., K_6C_{60}) is an insulator with the fullerenes in a bcc structure. The case $x=4$ is an insulator with the body-centred tetragonal structure, while the case $x=3$ is a metal with the fullerenes in an fcc structure. K_3C_{60} , Rb_3C_{60} , and Cs_3C_{60} are superconductors at low temperatures (see below *Low-temperature phenomena: Superconductivity*).

Alloys are solid mixtures of atoms with metallic properties. The definition includes both amorphous and crystalline solids. Although many pairs of elements will mix together as solids, many pairs will not. Almost all chemical entities can be mixed in liquid form. But cooling a

Alloys

Binary
crystal
structures

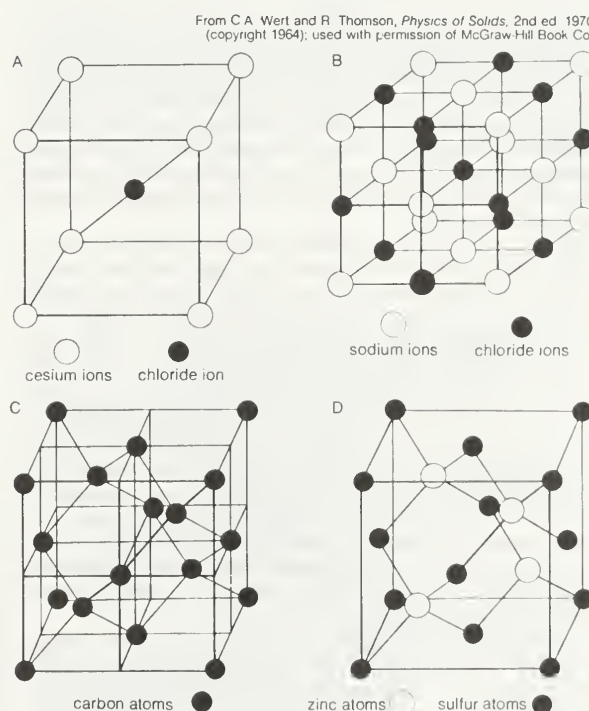


Figure 7: Crystal structures.

There is an equal number of the two types of ions in the unit cell of the (A) cesium chloride, (B) sodium chloride, and (D) zinc blende arrangements. The diamond arrangement is shown in (C). If both atoms are identical in (A), the structure is body-centred cubic.

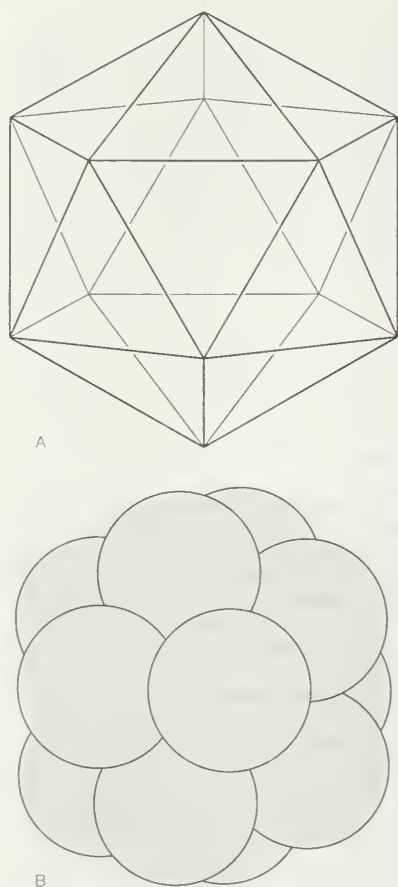


Figure 8: The icosahedral arrangement of boron molecules.

(A) Ions are located at the 12 vertices; the lines represent bonds between them. (B) The same molecule drawn with the atoms as spheres.

(B) From R.W.G. Wyckoff, *Crystal Structures*, 2nd ed., vol. 1, copyright © 1963 and 1965 John Wiley & Sons, Inc. reprinted by permission of John Wiley & Sons, Inc.

liquid to form a solid often results in phase separation; a polycrystalline material is obtained in which each grain is purely one atom or the other. Extremely rapid cooling can produce an amorphous alloy. Some pairs of elements form alloys that are metallic crystals. They have useful properties that differ from those exhibited by the pure elements. For example, alloying makes a metal stronger; for this reason alloys of gold, rather than the pure metal, are used in jewelry.

Atoms tend to form crystalline alloys when they are of similar size. The sizes of atoms are not easy to define, however, because atoms are not rigid objects with sharp boundaries. The outer part of an atom is composed of electrons in bound orbits; the average number of electrons decreases gradually with increasing distance from the nucleus. There is no point that can be assigned as the precise radius of the atom. Scientists have discovered, however, that each atom in a solid has a characteristic radius that determines its preferred separation from neighbouring atoms. For most types of atom this radius is constant, even in different solids. An empirical radius is assigned to each atom for bonding considerations, which leads to the concept of atomic size. Atoms readily make crystalline alloys when the radii of the two types of atoms agree to within roughly 15 percent.

Two kinds of ordering are found in crystalline alloys. Most alloys at low temperature are binary crystals with perfect ordering. An example is the alloy of copper and zinc. Copper is fcc, whereas zinc is hcp. A 50-percent-zinc–50-percent-copper alloy has a different structure— β -brass. At low temperatures it has the cesium chloride structure: a bcc lattice with alternating atoms of copper and zinc and a cubic unit cell. If the temperature is raised above 470° C, however, a phase transition to another crystalline state

occurs. The ordering at high temperature is also bcc, but now each site has equal probability of having a copper or zinc atom. The two types of atoms randomly occupy each site, but there is still long-range order. At all temperatures, thousands of atoms away from a site, the location of the atom site can be predicted with certainty. At temperatures below 470° C one also knows whether that site will be occupied by a copper or zinc atom, while above 470° C there is an equal likelihood of finding either atom. The high-temperature phase is crystalline but disordered. The disorder phase is obtained through a partial melting, not into a liquid state but into a less ordered one. This behaviour is typical of metal alloys. Other common alloys are steel, an alloy of iron and carbon; stainless steel, an alloy of iron, nickel (Ni), and chromium (Cr); pewter and solder, alloys of tin and lead (Pb); and britannia metal, an alloy of tin, antimony, and copper.

A crystal is never perfect; a variety of imperfections can mar the ordering. A defect is a small imperfection affecting a few atoms. The simplest type of defect is a missing atom and is called a vacancy. Since all atoms occupy space, extra atoms cannot be located at the lattice sites of other atoms, but they can be found between them; such atoms are called interstitials. Thermal vibrations may cause an atom to leave its original crystal site and move into a nearby interstitial site, creating a vacancy-interstitial pair. Vacancies and interstitials are the types of defects found in a pure crystal. In another defect, called an impurity, an atom is present that is different from the host crystal atoms. Impurities may either occupy interstitial spaces or substitute for a host atom in its lattice site.

There is no sharp distinction between an alloy and a crystal with many impurities. An alloy results when a sufficient number of impurities are added that are soluble in the host metal. However, most elements are not soluble in most crystals. Crystals generally can tolerate a few impurities per million host atoms. If too many impurities of the insoluble variety are added, they coalesce to form their own small crystallite. These inclusions are called precipitates and constitute a large defect.

Germanium is a common impurity in silicon. It prefers the same tetrahedral bonding as silicon and readily substitutes for silicon atoms. Similarly, silicon is a common impurity in germanium. No large crystal can be made without impurities; the purest large crystal ever grown was made of germanium. It had about 10^{10} impurities in each cubic centimetre of material, which is less than one impurity for each trillion atoms.

Impurities often make crystals more useful. In the absence of impurities, α -alumina is colourless. Iron and titanium impurities impart to it a blue colour, and the resulting gem-quality mineral is known as sapphire. Chromium impurities are responsible for the red colour characteristic of rubies, the other gem of α -alumina. Pure semiconductors rarely conduct electricity well at room temperatures. Their ability to conduct electricity is caused by impurities. Such impurities are deliberately added to silicon in the manufacture of integrated circuits. In fluorescent lamps the visible light is emitted by impurities in the phosphors (luminescent materials).

Other imperfections in crystals involve many atoms. Twinning is a special type of grain boundary defect, in which a crystal is joined to its mirror image. Another kind of imperfection is a dislocation, which is a line defect that may run the length of the crystal. One of the many types of dislocations is due to an extra plane of atoms that is inserted somewhere in the crystal structure. Another type, called an edge dislocation, is shown in Figure 9. This line defect occurs when there is a missing row of atoms. In the figure the crystal arrangement is perfect on the top and on the bottom. The defect is the row of atoms missing from region *b*. This mistake runs in a line that is perpendicular to the page and places a strain on region *a*.

Dislocations are formed when a crystal is grown, and great care must be taken to produce a crystal free of them. Dislocations are stable and will exist for years. They relieve mechanical stress. If one presses on a crystal, it will accommodate the induced stress by growing dislocations at the surface, which gradually move inward. Dislocations

Crystal defects

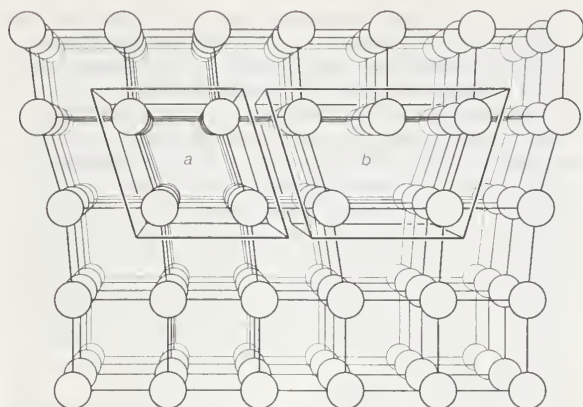


Figure 9: Crystalline lattice defect.

An edge dislocation occurs when there is a missing row of atoms as shown in region *b*. Region *a* is strained.

From J. Ziman, "The Thermal Properties of Materials," copyright © 1967 by Scientific American, Inc., all rights reserved.

make a crystal mechanically harder. When a metal bar is cold-worked by rolling or hammering, dislocations and grain boundaries are introduced; this causes the hardening.

Determination
of crystal
structures

Crystal structures are determined by scattering experiments using a portion of the crystal as the target. A beam of particles is sent toward the target, and upon impact some of the particles scatter from the crystal and ricochet in various directions. A measurement of the scattered particles provides raw data, which is then computer-processed to give a picture of the atomic arrangements. The positions are then inferred from the computer-analyzed data.

Max von Laue first suggested in 1912 that this measurement could be done using X rays, which are electromagnetic radiation of very high frequency. High frequencies are needed because these waves have a short wavelength. Von Laue realized that atoms have a spacing of only a few angstroms (1 angstrom [Å] is 10^{-10} metre, or 3.94×10^{-9} inch). In order to measure atomic arrangements, the particles scattering from the target must also have a wavelength of a few angstroms. X rays are required when the beam consists of electromagnetic radiation. The X rays only scatter in certain directions, and there are many X rays associated with each direction. The scattered particles appear in spots corresponding to locations where the scattering from each identical atom produces an outgoing wave that has all the wavelengths in phase. Figure 10 shows incoming waves in phase. The scattering from atom A_2 has a longer path than that from atom A_1 . If this additional path has a length ($AB + BC$) that is an exact multiple of the wavelength, then the two outgoing waves are in phase and reinforce each other. If the scattering angle is changed slightly, the waves no longer add coherently and begin to cancel one another. Combining the scattered radiation from all the atoms in the crystal causes all the outgoing waves to add coherently in certain directions and produce a strong signal in the scattered wave. If the extra path length ($AB + BC$) is five wavelengths, for example, the spot appears in one place. If it is six wavelengths, the spot is elsewhere. Thus, the different spots correspond to the different multiples of the wavelength of the X ray. The measurement produces two types of information: the directions of the spots and their intensity. This information is insufficient to deduce the exact crystal structure, however, as there is no algorithm by which the computer can go directly from the data to the structure. The crystallographer must propose various structures and compute how they would scatter the X rays. The theoretical results are compared with the measured one, and the theoretical arrangement is chosen that best fits the data. Although this procedure is fast when there are only a few atoms in a unit cell, it may take months or years for complex structures. Some protein molecules, for instance, have hundreds of atoms. Crystals of the proteins are grown, and X rays are used to measure the structure. The goal is to determine how the atoms are arranged in the protein, rather than how the proteins are arranged in the crystal.

Beams of neutrons may also be used to measure crystal

structure. The beam of neutrons is obtained by drilling a hole in the side of a nuclear reactor. The energetic neutrons created in nuclear fission escape through the hole. The motion of elementary particles is governed by quantum, or wave, mechanics. Each neutron has a wavelength that depends on its momentum. The scattering directions are determined by the wavelength, as is the case with X rays. The wavelengths for neutrons from a reactor are suitable for measuring crystal structures.

X rays and neutrons provide the basis for two competing technologies in crystallography. Although they are similar in principle, the two methods have some differences. X rays scatter from the electrons in the atoms so that more electrons result in more scattering. X rays easily detect atoms of high atomic number, which have many electrons, but cannot readily locate atoms with few electrons. In hydrogen-bonded crystals, X rays do not detect the protons at all. Neutrons, on the other hand, scatter from the atomic nucleus. They scatter readily from protons and are excellent for determining the structure of hydrogen-bonded solids. One drawback to this method is that some nuclei absorb neutrons completely, and there is little scattering from these targets.

Beams of electrons can also be used to measure crystal structure, because energetic electrons have a wavelength that is suitable for such measurements. The problem with electrons is that they scatter strongly from atoms. Proper interpretation of the experimental results requires that an electron scatter only from one atom and leave the crystal without scattering again. Low-energy electrons scatter many times, and the interpretation must reflect this. Low-energy electron diffraction (LEED) is a technique in which a beam of electrons is directed toward the surface. The scattered electrons that reflect backward from the surface are measured. They scatter many times before leaving backward but mainly leave in a few directions that appear as "spots" in the measurements. An analysis of the varied spots gives information on the crystalline arrangement. Because the electrons are scattered strongly by the atoms in the first few layers of the surface, the measurement gives only the arrangements of atoms in these layers. It is assumed that the same structure is repeated throughout the crystal. Another scattering experiment involves electrons of extremely high energy. The scattering rate decreases as the energy of the electron increases, so that very energetic electrons usually scatter only once. Various electron microscopes are constructed on this principle. A photograph taken by an electron microscope is shown in Figure 11. The bright spots are individual atoms. The dark region on the bottom is a crystal of silicon, and the brighter regions on top are epitaxially grown films of an alloy of germanium and silicon. The arrow points to dislocations formed because of the lattice mismatch.

Electron-diffraction
techniques

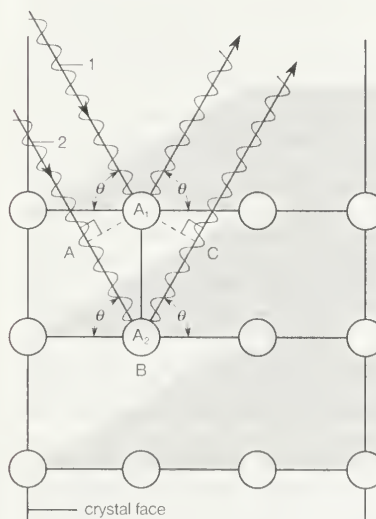


Figure 10: Incident rays (1 and 2) at angle θ on the planes of atoms in a crystal. Rays reinforce if their difference in path length ($AB + BC$) is an integer times the wavelength of the X ray.

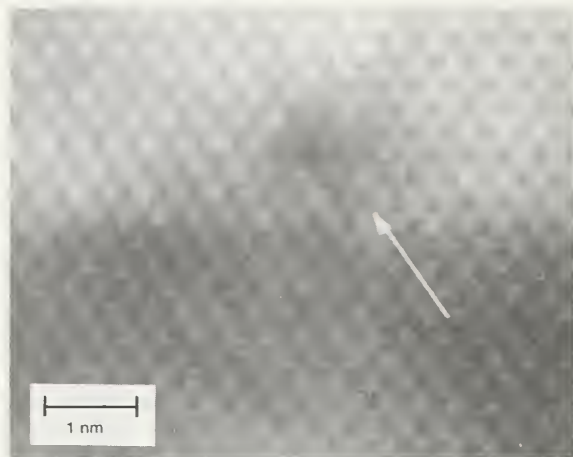


Figure 11: An electron-microscope image. The bright spots are individual atoms. The brighter upper regions are an epitaxially grown film of a germanium-silicon alloy; the darker region below them is the silicon substrate. The arrow points to a dislocation that resulted from lattice mismatch.

By courtesy of S.J. Pennycook and M.F. Chisholm, Oak Ridge National Laboratory.

TYPES OF BONDS

The properties of a solid can usually be predicted from the valence and bonding preferences of its constituent atoms. Four main bonding types are discussed here: ionic, covalent, metallic, and molecular. Hydrogen-bonded solids, such as ice, make up another category that is important in a few crystals. There are many examples of solids that have a single bonding type, while other solids have a mixture of types, such as covalent and metallic or covalent and ionic.

Ionic bonds. Sodium chloride exhibits ionic bonding. The sodium atom has a single electron in its outermost shell, while chlorine needs one electron to fill its outer shell. Sodium donates one electron to chlorine, forming a sodium ion (Na^+) and a chlorine ion (Cl^-). Each ion thus attains a closed outer shell of electrons and takes on a spherical shape. In addition to having filled shells and a spherical shape, the ions of an ionic solid have integer valence. An ion with positive valence is called a cation. In an ionic solid the cations are surrounded by ions with negative valence, called anions. Similarly, each anion is surrounded by cations. Since opposite charges attract, the preferred bonding occurs when each ion has as many neighbours as possible, consistent with the ion radii. Six or eight nearest neighbours are typical; the number depends on the size of the ions and not on the bond angles. The alkali halide crystals are binaries of the AH type, where A is an alkali ion (lithium [Li], sodium, potassium, rubidium, or cesium) and H is a halide ion (fluorine, chlorine, bromine, or iodine). The crystals have ionic bonding, and each ion has six or eight neighbours. Metal ions in the alkaline earth series (magnesium [Mg], calcium [Ca], barium [Ba], and strontium [Sr]) have two electrons in their outer shells and form divalent cations in ionic crystals. The chalcogenides (oxygen, sulfur, selenium, and tellurium) need two electrons to fill their outer p -shell. (Electron shells are divided into subshells, designated as s , p , d , f , g , and so forth. Each subshell is divided further into orbitals.) Two electrons are transferred from the cations to the anions, leaving each with a closed shell. The alkaline earth chalcogenides form ionic binary crystals such as barium oxide (BaO), calcium sulfide (CaS), barium selenide (BaSe), or strontium oxide (SrO). They have the same structure as sodium chloride, with each atom having six neighbours. Oxygen can be combined with various cations to form a large number of ionically bonded solids.

Covalent bonds. Silicon, carbon, germanium, and a few other elements form covalently bonded solids. In these elements there are four electrons in the outer sp -shell, which is half filled. (The sp -shell is a hybrid formed from one s and one p subshell.) In the covalent bond an atom shares one valence (outer-shell) electron with each of its four nearest neighbour atoms. The bonds are highly directional

and prefer a tetrahedral arrangement. A covalent bond is formed by two electrons—one from each atom—located in orbitals between the ions. Insulators, in contrast, have all their electrons within shells inside the atoms.

The perpetual spin of an electron is an important aspect of the covalent bond. From a vantage point above the spinning particle, counterclockwise rotation is designated spin-up, while clockwise rotation is spin-down. A fundamental law of quantum physics is the Pauli exclusion principle, which states that no two electrons can occupy the same point in space at the same time with the same direction of spin. In a covalent bond two electrons occupy the same small volume of space (*i.e.*, the same orbital) at all times, so they must have opposite spin: one up and one down. The exclusion principle is then satisfied, and the resulting bond is strong.

In graphite the carbon atoms are arranged in parallel sheets, and each atom has only three near neighbours. The covalent bonds between adjacent carbons within each layer are quite strong and are called σ bonds. The fourth valence electron in carbon has its orbital perpendicular to the plane. This orbital bonds weakly with the similar orbitals on all three neighbours, forming π bonds. The four bonds for each carbon atom in the graphite structure are not arranged in a tetrahedron; three are in a plane. The planar arrangement results in strong bonding, although not as strong as the bonding in the diamond configuration. The bonding between layers is quite weak and arises from the van der Waals interaction; there is much slippage parallel to the layers. Diamond and graphite form an interesting contrast: diamond is the hardest material in nature and is used as an abrasive, while graphite is used as a lubricant.

Besides the elemental semiconductors, such as silicon and germanium, some binary crystals are covalently bonded. Gallium has three electrons in the outer shell, while arsenic lacks three. Gallium arsenide (GaAs) could be formed as an insulator by transferring three electrons from gallium to arsenic; however, this does not occur. Instead, the bonding is more covalent, and gallium arsenide is a covalent semiconductor. The outer shells of the gallium atoms contribute three electrons, and those of the arsenic atoms contribute five, providing the eight electrons needed for four covalent bonds. The centres of the bonds are not at the midpoint between the ions but are shifted slightly toward the arsenic. Such bonding is typical of the III-V semiconductors—*i.e.*, those consisting of one element from the third column of the periodic table and one from the fifth column. Elements from the third column (boron, aluminum, gallium, and indium) contribute three electrons, while the fifth-column elements (nitrogen, phosphorus, arsenic, and antimony) contribute five electrons. All III-V semiconductors are covalently bonded and typically have the zinc blende structure with four neighbours per atom. Most common semiconductors favour this arrangement.

The factor that determines whether a binary crystal will act as an insulator or a semiconductor is the valence of its constituent atoms. Ions that donate or accept one or two valence electrons form insulators. Those that have three to five valence electrons tend to have covalent bonds and form semiconductors. There are exceptions to these rules, however, as is the case with the IV-VI semiconductors such as lead sulfide. Heavier elements from the fourth column of the periodic table (germanium, tin, and lead) combine with the chalcogenides from the sixth row to form good binary semiconductors such as germanium telluride (GeTe) or tin sulfide (SnS). They have the sodium chloride structure, where each atom has six neighbours. Although not tetrahedrally bonded, they are good semiconductors.

Filled atomic shells with d -orbitals have an important role in covalent bonding. Electrons in atomic orbits have angular momentum (L), which is quantized in integer (n) multiples of Planck's constant h : $L = nh$. Electron orbitals with $n = 0$ are called s -states, with $n = 1$ are p -states, and with $n = 2$ are d -states. Silver and copper ions have one valence electron outside their closed shells. The outermost filled shell is a d -state and affects the bonding. Eight binary crystals are formed from the copper and silver

Electron spin

Bonding in binary semiconductors

halides. Three (AgF, AgCl, AgBr) have the sodium chloride structure with six neighbours. The other five (AgI, CuF, CuCl, CuBr, CuI) have the zinc blende structure with four neighbours. The bonding in this group of solids is on the borderline between covalent and ionic, since the crystals prefer both types of bonds. The alkali metal halides exhibit somewhat different behaviour. The alkali metals are also monovalent cations, but their halides are strictly ionic. The difference in bonding between the alkali metals on the one hand and silver and copper on the other hand is that silver and copper have filled *d*-shells while the alkalis have filled *p*-shells. Since the *d*-shells are filled, they do not covalently bond. This group of electrons is, however, highly polarizable, which influences the bonding of the valence electrons. Similar behaviour is found for zinc and cadmium, which have two valence electrons outside a filled *d*-shell. They form binary crystals with the chalcogenides, which have tetrahedral bonding. In this case the covalent bonding seems to be preferred over the ionic bond. In contrast, the alkaline earth chalcogenides, which are also divalent, have outer *p*-shells and are ionic. The zinc and cadmium chalcogenides are covalent, as the outer *d*-shell electrons of the two cations favour covalent bonding.

Metallic bonds. Metallic bonds fall into two categories. The first is the case in which the valence electrons are from the *sp*-shells of the metal ions; this bonding is quite weak. In the second category the valence electrons are from partially filled *d*-shells, and this bonding is quite strong. The *d*-bonds dominate when both types of bonding are present.

The simple metals are bonded with *sp*-electrons. The electrons of these metal atoms are in filled atomic shells except for a few electrons that are in unfilled *sp*-shells. The electrons from the unfilled shells are detached from the metal ion and are free to wander throughout the crystal. They are called conduction electrons, since they are responsible for the electrical conductivity of metals. Although the conduction electrons may roam anywhere in the crystal, they are distributed uniformly throughout the entire solid. Any large imbalance of charge is prevented by the strong electrical attraction between the negative electrons and the positive ions, plus the strong repulsion between electrons. The phrase electron correlation describes the correlated movements of the electrons; the motion of each electron depends on the positions of neighbouring electrons. Electrons have strong short-range order with one another. Correlation ensures that each unit cell in the crystal has, on the average, the number of electrons needed to cancel the positive charge of the cation so that the unit cell is electrically neutral.

Cohesive energy is the energy gained by arranging the atoms in a crystalline state, as compared with the gas state. Insulators and semiconductors have large cohesive energies; these solids are bound together strongly and have good mechanical strength. Metals with electrons in *sp*-bonds have very small cohesive energies. This type of metallic bond is weak; the crystals are barely held together. Single crystals of simple metals such as sodium are mechanically weak. At room temperature the crystals have the mechanical consistency of warm butter. Special care must be used in handling these crystals, because they are easily distorted. Metals such as magnesium or aluminum must be alloyed or polycrystalline to have any mechanical strength. Although the simple metals are found in a variety of structures, most are in one of the three closest-packed structures: fcc, bcc, and hcp. Theoretical calculations show that the cohesive energy of a given metal is almost the same in each of the different crystal arrangements; therefore, crystal arrangements are unimportant in metals bound with electrons from *sp*-shells.

A different type of metallic bonding is found in transition metals, which are metals whose atoms are characterized by unfilled *d*-shells. The *d*-orbitals are more tightly bound to an ion than the *sp*-orbitals. Electrons in *d*-shells do not wander away from the ion. The *d*-orbitals form a covalent bond with the *d*-orbitals on the neighbouring atoms. The bonding of *d*-orbitals does not occur in a tetrahedral arrangement but has a different directional preference. In metals the bonds from *d*-orbitals are not completely filled

with electrons. This situation is different from the tetrahedral bonds in semiconductors, which are filled with eight electrons. In transition metals the covalent bonds formed with the *d*-electrons are much stronger than the weak bonds made with the *sp*-electrons of simple metals. The cohesive energy is much larger in transition metals. Titanium, iron, and tungsten, for example, have exceptional mechanical strength. Crystal arrangements are important in the behaviour of the transition metals and occur in the close-packed fcc, bcc, or hcp arrangements.

Molecular binding. The Dutch physicist Johannes D. van der Waals first proposed the force that binds molecular solids. Any two atoms or molecules have a force of attraction (*F*) that varies according to the inverse seventh power of the distance *R* between the centres of the atoms or molecules: $F = -C/R^7$, where *C* is a constant. The force, known as the van der Waals force, declines rapidly with the distance *R* and is quite weak. If the atoms or molecules have a net charge, there is a strong force whose strength varies according to Coulomb's law as the inverse second power of the separation distance: $F = -C'/R^2$, where *C'* is a constant. This force provides the binding in ionic crystals and some of the binding in metals. Coulomb's law does not apply to atoms or molecules without a net charge. Molecules with a dipole moment, such as water, have a strong attractive force owing to the interactions between the dipoles. For atoms and molecules with neither net charges nor dipole moments, the van der Waals force provides the crystal binding. The force of gravity also acts between neutral atoms and molecules, but it is far too weak to bind molecules into crystals.

The van der Waals force is caused by quantum fluctuations. Two neighbouring atoms that are each fluctuating can lower their joint energy by correlating their fluctuations. The van der Waals force arises from correlations in their dipole fluctuations. Electrons bound in atomic orbits are in constant motion around the nucleus, and the distribution of charges in the atom changes constantly as the electrons move, owing to quantum fluctuations. One fluctuation might produce a momentary electric dipole moment (*i.e.*, a separation of charges) on an atom if a majority of its electrons are on one side of the nucleus. The dipole moment creates an electric field on a neighbouring atom; this field will induce a dipole moment on the second atom. The two dipoles attract one another via the van der Waals interaction. Since the force depends on the inverse seventh power of the distance, it declines rapidly with increasing distance. Atoms have a typical radius of one to three angstroms. The van der Waals force binds atoms and molecules within a few angstroms of each other; beyond that distance the force is negligible. Although weak, the van der Waals force is always present and is important in cases where the other forces are absent.

Hydrogen is rarely found as a single atom. Instead it forms diatomic molecules (H_2), which are gaseous at room temperature. At lower temperatures the hydrogen becomes a liquid and at about 20 K turns into a solid. The molecule retains its identity in the liquid and solid states. The solid exists as a molecular crystal of covalently bound H_2 molecules. The molecules attract one another by van der Waals forces, which provide the crystal binding. Helium, the second element in the periodic table, has two electrons, which constitute a filled atomic shell. In its liquid and solid states, the helium atoms are bound together by van der Waals forces. In fact, all the rare gases (helium, neon, argon, krypton, and xenon) are molecular crystals with the binding provided by van der Waals forces.

Many organic molecules form crystals where the molecules are bound by van der Waals forces. In methane (CH_4), a central carbon makes a covalent bond with each hydrogen atom, forming a tetrahedron. In crystalline methane the molecules are arranged in the fcc structure. Benzene (C_6H_6) has the carbon atoms in a hexagonal ring; each carbon has three coplanar σ bonds, as in graphite, where two bonds are with neighbouring carbon atoms and the third bond is with a hydrogen atom. Crystalline benzene has four molecules per unit cell in a complex arrangement. Fullerene and the rare gas atoms are spherical, and the crystalline arrangement corresponds to the closest

The van der Waals force

Electron correlation

Bonding in crystals of organic molecules

packing of spheres. Most organic molecules, however, are not spherical and display irregular shapes. For odd-shaped molecules, the van der Waals interaction depends on the rotational orientation of the two molecules. In order to maximize the force, the molecules in the crystal have unusual arrangements, as in the case of benzene.

Hydrogen bonding. Hydrogen bonding is important in a few crystals, notably in ice. With its lone electron, a hydrogen atom usually forms a single covalent bond with an electronegative atom. In the hydrogen bond the atom is ionized to a proton. The proton sits between two anions and joins them. Hydrogen bonding occurs with only the most electronegative ions: nitrogen, oxygen, and fluorine. In water the hydrogen links pairs of oxygen ions. Water is found in many different crystal structures, but they all have the feature that the hydrogen atoms sit between pairs of oxygen. Another hydrogen-bonded solid is hydrogen fluoride (HF), in which the hydrogen atom (proton) links pairs of fluorines.

CRYSTAL GROWTH

The earliest crystal grower was nature. Many excellent crystals of minerals formed in the geologic past are found in mines and caves throughout the world. Most precious and semiprecious stones are well-formed crystals. Early efforts to produce synthetic crystals were concentrated on making gems. Synthetic ruby was grown by the French scientist Marc Antoine Augustin Gaudin in 1873. Since about 1950 scientists have learned to grow in the laboratory crystals of quality equal or superior to those found in nature. New techniques for growth are continually being developed, and crystals with three or more atoms per unit cell are continually being discovered.

Vapour growth. Crystals can be grown from a vapour when the molecules of the gas attach themselves to a surface and move into the crystal arrangement. Several important conditions must be met for this to occur. At constant temperature and equilibrium conditions, the average number of molecules in the gas and solid states is constant: molecules leave the gas and attach to the surface at the same rate that they leave the surface to become gas molecules. For crystals to grow, the gas-solid chemical system must be in a nonequilibrium state such that there are too many gaseous molecules for the conditions of pressure and temperature. This state is called supersaturation. Molecules are more prone to leave the gas than to rejoin it, so they become deposited on the surface of the container. Supersaturation can be induced by maintaining the crystal at a lower temperature than the gas. A critical stage in the growth of a crystal is seeding, in which a small piece of crystal of the proper structure and orientation, called a seed, is introduced into the container. The gas molecules find the seed a more favourable surface than the walls and preferentially deposit there. Once the molecule is on the surface of the seed, it wanders around this surface to find the preferred site for attachment. Growth proceeds one molecule at a time and one layer at a time. The process is slow: it takes days to grow a small crystal. Crystals are grown at temperatures well below the melting point to reduce the density of defects. The advantage of vapour growth is that very pure crystals can be grown by this method, while the disadvantage is that it is slow.

Most clouds in the atmosphere are ice crystals that form by vapour growth from water molecules. Most raindrops are crystals as they begin descending but thaw during their fall to Earth. Seeding for rain—accomplished by dropping silver iodide crystals from airplanes—is known to induce precipitation. In the laboratory, vapour growth is usually accomplished by flowing a supersaturated gas over a seed crystal. Quite often a chemical reaction at the surface is needed to deposit the atoms. Crystals of silicon can be grown by flowing chlorosilane (SiCl_4) and hydrogen (H_2) over a seed crystal of silicon. Hydrogen acts as the buffer gas by controlling the temperature and rate of flow. The molecules dissociate on the surface in a chemical reaction that forms hydrogen chloride (HCl) molecules. Hydrogen chloride molecules leave the surface, while silicon atoms remain to grow into a crystal. Binary crystals such as gallium arsenide (GaAs) are grown by a similar method. One

process employs gallium chloride (GaCl_3) as the gallium carrier. Arsenic is provided by molecules such as arsenous chloride (AsCl_3), arsine (AsH_3), or As_4 (yellow arsenic). These molecules, with hydrogen as the buffer gas, grow crystals of gallium arsenide while forming gas molecules such as gallium trichloride (GaCl_3) and hydrogen chloride. Trimethylgallium, $(\text{CH}_3)_3\text{Ga}$, is another molecule that can be used to deliver gallium to the surface.

Growth from solution. Large single crystals may be grown from solution. In this technique the seed crystal is immersed in a solvent that contains typically about 10–30 percent of the desired solute. The choice of solvent usually depends on the solubility of the solute. The temperature and pH (a measure of acidity or basicity) of the solution must be well controlled. The method is faster than vapour growth, because there is a higher concentration of molecules at the surface in a liquid as compared to a gas, but it is still relatively slow.

Growth from the melt. This method is the most basic. A gas is cooled until it becomes a liquid, which is then cooled further until it becomes a solid. Polycrystalline solids are typically produced by this method unless special techniques are employed. In any case, the temperature must be controlled carefully. Large crystals can be grown rapidly from the liquid elements using a popular method invented in 1918 by the Polish scientist Jan Czochralski and called crystal pulling. One attaches a seed crystal to the bottom of a vertical arm such that the seed is barely in contact with the material at the surface of the melt. A modern Czochralski apparatus is shown in Figure 12A. The arm is raised slowly, and a crystal grows underneath at the interface between the crystal and the melt. Usually the crystal is rotated slowly, so that inhomogeneities in the liquid are not replicated in the crystal. Large-diameter crystals of silicon are grown in this way for use as computer chips. Based on measurements of the weight of the crystal during the pulling process, computer-controlled apparatuses can vary the pulling rate to produce any de-

Crystal pulling

Seeding

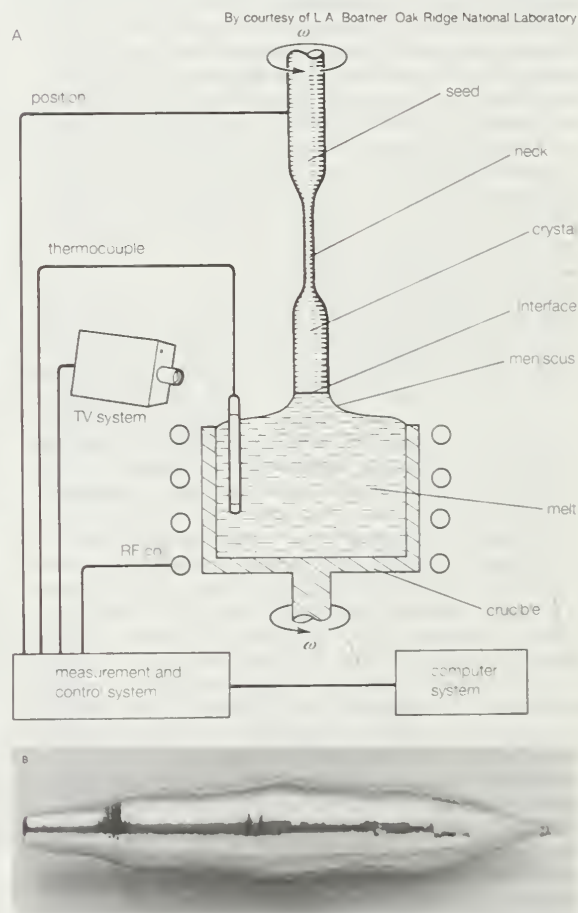


Figure 12: Crystal pulling using the Czochralski method. (A) A schematic view of a modern apparatus. (B) A crystal of stainless steel grown by this method.

sired diameter. Crystal pulling is the least expensive way to grow large amounts of pure crystal. A photograph of a single crystal of stainless steel grown by the Czochralski method is shown in Figure 12B. The original seed is on the right tip. Binary crystals can also be pulled; for example, synthetic sapphire crystals can be pulled from molten alumina. Special care is required to grow binary and other multicomponent crystals; the temperature must be precisely controlled because such crystals may be grown only at a single, extremely high temperature. The melt has a tendency to be inhomogeneous, since the two liquids may try to separate by gravity.

The Bridgman method (named after the American scientist Percy Williams Bridgman) is also widely used for growing large single crystals. The molten material is put into a crucible, often of silica, which has a cylindrical shape with a conical lower end. Heaters maintain the molten state. As the crucible is slowly lowered into a cooler region, a crystal starts growing in the conical tip. The crucible is lowered at a rate that matches the growth of the crystal, so that the interface between crystal and melt is always at the same temperature. The rate of moving the crucible depends on the temperature and the material. When done successfully, the entire molten material in the crucible grows into a single large crystal. One disadvantage of the method is that excess impurities are pushed out of the crystal during growth. A layer of impurities grows at the interface between melt and solid as this surface moves up the melt, and the impurities become concentrated in the higher part of the crystal.

Epitaxy Epitaxy is the technique of growing a crystal, layer by layer, on the atomically flat surface of another crystal. In homoepitaxy a crystal is grown on a substrate of the same material. Silicon layers of different impurity content, for example, are grown on silicon substrates in the manufacture of computer chips. Heteroepitaxy, on the other hand, is the growth of one crystal on the substrate of another. Silicon substrates are often used since they are readily available in atomically smooth form. Many different semiconductor crystals can be grown on silicon, such as gallium arsenide, germanium, cadmium telluride (CdTe), and lead telluride (PbTe). Any flat substrate can be used for epitaxy, however, and insulators such as rock salt (NaCl) and magnesium oxide (MgO) are also used.

Molecular-beam epitaxy, commonly abbreviated as MBE, is a form of vapour growth. The field began when the American scientist John Read Arthur reported in 1968 that gallium arsenide could be grown by sending a beam of gallium atoms and arsenic molecules toward the flat surface of a crystal of the molecule. The amount of gas molecules can be controlled to grow just one layer, or just two, or any desired amount. This method is slow, since molecular beams have low densities of atoms. Chemical vapour deposition (CVD) is another form of epitaxy that makes use of the vapour growth technique. Also known as vapour-phase epitaxy (VPE), it is much faster than MBE since the atoms are delivered in a flowing gas rather than in a molecular beam. Synthetic diamonds are grown by CVD. Rapid growth occurs when methane (CH_4) is mixed with atomic hydrogen gas, which serves as a catalyst. Methane dissociates on a heated surface of diamond. The carbon remains on the surface, and the hydrogen leaves as a molecule. Growth rates are several micrometres (1 micrometre is equal to 0.00004 inch) per hour. At that rate, a stone 1 centimetre (0.4 inch) thick is grown in 18 weeks. CVD diamonds are of poor quality as gemstones but are important electronic materials. Because hydrogen is found in nature as a molecule rather than as a single atom, making atomic hydrogen gas is the major expense in growing CVD diamonds. Liquid-phase epitaxy (LPE) uses the solution method to grow crystals on a substrate. The substrate is placed in a solution with a saturated concentration of solute. This technique is used to grow many crystals employed in modern electronics and optoelectronic devices, such as gallium arsenide, gallium aluminum arsenide, and gallium phosphide.

An important concern in successful epitaxy is matching lattice distances. If the spacing between atoms in the substrate is close to that of the top crystal, then that crystal

will grow well; a small difference in lattice distance can be accommodated as the top crystal grows. When the lattice distances are different, however, the top crystal becomes deformed, since structural defects such as dislocations appear (see Figure 9). Although few crystals share the same lattice distance, a number of examples are known. Aluminum arsenide and gallium arsenide have the same crystal structure and the same lattice parameters to within 0.1 percent; they grow excellent crystals on one another. Such materials, known as superlattices, have a repeated structure of n layers of GaAs, m layers of AlAs, n layers of GaAs, m layers of AlAs, and so forth. Superlattices represent artificially created structures that are thermodynamically stable; they have many applications in the modern electronics industry. Another lattice-matched epitaxial system is mercury telluride (HgTe) and cadmium telluride (CdTe). These two semiconductors form a continuous semiconductor alloy $\text{Cd}_x\text{Hg}_{1-x}\text{Te}$, where x is any number between 0 and 1. This alloy is used as a detector of infrared radiation and is incorporated in particular in night-vision goggles.

Superlattices

Dendritic growth. At slow rates of crystal growth, the interface between melt and solid remains planar, and growth occurs uniformly across the surface. At faster rates of crystal growth, instabilities are more likely to occur; this leads to dendritic growth. Solidification releases excess energy in the form of heat at the interface between solid and melt. At slow growth rates, the heat leaves the surface by diffusion. Rapid growth creates more heat, which is dissipated by convection (liquid flow) when diffusion is too slow. Convection breaks the planar symmetry so that crystal growth develops along columns, or "fingers," rather than along planes. Each crystal has certain directions in which growth is fastest, and dendrites grow in these directions. As the columns grow larger, their surfaces become flatter and more unstable. Columns start growing from other columns in the pattern shown in Figure 13. This feather or tree structure is characteristic of dendritic growth. Snowflakes are an example of crystals that result from dendritic growth.

ELECTRIC PROPERTIES

The German physicist Georg Simon Ohm discovered the basic law of electric conduction, which is now called Ohm's law. His law relates the voltage (V , measured in volts), the current (I , in amperes), and the resistance (R , in ohms) according to the formula $V = RI$. A current

From A.H. Cottrell, *An Introduction of Metallurgy*, © 1967 Edward Arnold publishers

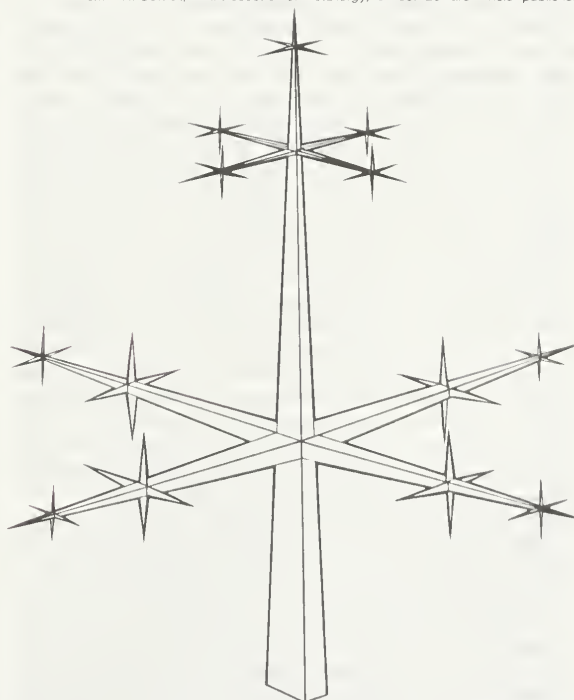


Figure 13: A possible pattern during dendritic growth of crystals.

I through a solid induces a voltage V ; the resistance R is the constant of proportionality. The value of R is an important factor in the design of electrical circuits. It is determined by the shape of the resistor: a long narrow object has more resistance than a short wide one of the same material. For solids, the important parameter is the resistivity ρ , which is given in units of ohm-metres. It is the resistivity per volume unit and is independent of shape. The relationship between R and ρ is $R = \rho L/A$, where A is the area of the resistor and L is the length. These dimensions are measured in the direction of the current: L is the length of the current path, and A is the cross-sectional area. The resistance of a copper bar depends on its shape, but at a given temperature every piece of pure copper has the same resistivity. Thus the resistivity is a fundamental parameter of a material and is investigated by scientists. Resistivities of solids span a wide range of values. Certain metals have zero resistivity at low temperatures; they are called superconductors. At the other extreme, very good insulators such as sulfur and polystyrene have resistivities larger than one quadrillion ohm-metres. At room temperature, the metal with the lowest value of resistivity is silver, with $\rho = 1.6 \times 10^{-8}$ ohm-metre; the second best conductor is copper, with $\rho = 1.7 \times 10^{-8}$ ohm-metre. Copper, rather than silver, is used in household wires because of the high cost of silver.

Electrical conductivity σ is the inverse of resistivity and is measured in units of ohm-metre⁻¹. Electrical current is produced by the motion of charges. In crystals, electrical current is due to the motion of both ions and electrons. Ions move by hopping occasionally from site to site; all solids can conduct electricity in this manner. When the voltage is zero, there is no net current because the ions hop randomly in all directions. The imposition of a small voltage causes the ions to slightly favour one direction of motion, which leads to a net flow of charge in that direction; this constitutes an electrical current. The electricity conducted by this process is quite small and is usually negligible compared with that carried by the electrons. When an ion hops, it must migrate to a vacant site, which could be either an interstitial or a vacancy. Ionic conductivity can occur because hopping ions cause vacancies to move through the solid. An ion hops to the vacant site, thereby filling the vacancy, while creating a new one at the ion's former site. Repeating this process causes the position of the vacancy to migrate through the crystal. The motion of the vacancy arises from the motion of ions, which carry charge and contribute to electrical conductivity.

Ion hops are induced by thermal fluctuations. Most of the ions move within their lattice site, vibrating around this point. Temperature is defined as the average energy of this vibrational motion; the more the ions move, the higher the temperature. An individual ion at times moves slowly and at times vibrates quite rapidly but usually has an energy near the average value. Each ion shares its vibrational energy with its neighbouring ions. An ion typically has some neighbours with small vibrations and others with large ones. The average energy shared with the neighbours is close to the average energy of all the atoms. As a random process, however, it occasionally happens that all neighbours of an ion may have large vibrations, in which case the ion will acquire an unusually high energy. This energy may be high enough to cause it to leave its site and hop to a neighbouring site. A thermal fluctuation is the rare process in which the energy at a local site may be much higher or lower than the average energy in the crystal. Probability theory shows that the higher the temperature, the more frequent are these thermal fluctuations. Ions therefore hop more often at high temperature.

A few solids conduct electricity better by ion motion than by electron motion. These unusual materials are technologically important in making batteries. All batteries have two electrodes separated by an electrolyte, which is a material that conducts ions better than electrons. An example of a crystal electrolyte is β -alumina, which readily conducts monovalent cations such as silver (Ag^+) and sodium (Na^+). Among all ions, silver has the largest value of ionic conductivity in many different electronic insulators. The copper ion (Cu^+) forms the same type of

chemical bonds as does the silver ion, but the copper ion, because of its smaller radius, does not migrate as well within an electrolyte. Silver ions fit perfectly into the interstitial sites of the crystal lattices of several electrolytes, while the smaller copper ions permit the neighbouring ions to collapse around them, inhibiting further hopping. There are a few good conductors of the inexpensive copper ion that can be used as solid electrolytes in batteries. Silver is too costly and heavy to use in large-volume batteries such as those found in automobiles, but it is used in the smaller batteries that power devices such as hearing aids.

Electrons carry the basic unit of charge e , equal to 1.6022×10^{-19} coulomb. They have a small mass and move rapidly. Most electrons in solids are bound to the atoms in local orbits, but a small fraction of the electrons are available to move easily through the entire crystal. These so-called conduction electrons carry the electrical current. Solids with many conduction electrons are metals, while those with a few are semimetals or semiconductors. In insulators, nearly all the electrons are bound, and very few electrons are capable of carrying current. A typical metal has one or more conduction electrons in each atomic unit cell, a semiconductor may have only one conduction electron for each thousand unit cells, and an insulator may have one conduction electron per one million or one trillion unit cells.

The bonding properties of the individual atoms of a solid determine the behaviour of the bulk solid. The electrical properties of a solid can usually be predicted from the valence and bonding preferences of its atoms. In the argon atom, for example, all atomic shells are filled with electrons. The electrons of solid argon remain in the atomic shells; none are conduction electrons, and the electrical resistivity is therefore high. Solid argon, like all the rare gas solids, is a good insulator. A few conduction electrons are contributed by impurities, and so the conductivity, though small, is not zero. These conduction electrons move quite readily through the solid. The term mobility is used to describe how well a conduction electron moves through the solid in response to a voltage. Conductivity is the product of mobility, the electrical charge e , and the number N of conduction electrons per unit volume: $\sigma = Ne\mu$, where σ is the conductivity and μ is the mobility. The mobility of the rare gas solids is high, but their conductivity is nonetheless low because there is a small number of conduction electrons.

Like the rare gas solids, most ionic solids are electrical insulators. In sodium chloride, for example, each sodium atom donates its single valence electron to a chlorine atom, thus forming a solid composed of Na^+ and Cl^- ions. All electrons are in filled shells at low temperature, and in a perfect crystal there are no conduction electrons. Sodium chloride is thus an insulator with a very high resistivity. Some conduction electrons are provided by impurities or thermal excitations. At high temperatures large ion vibrations from thermal fluctuations may knock an electron out of a filled shell, upon which it becomes a conduction electron and contributes to the conductivity. The number of conduction electrons created by thermal excitations is small for most insulators. Although defects can be responsible for producing conduction electrons, they can also destroy the conducting ability of electrons by trapping them. The defects have local orbitals that provide a lower energy state for the electron than the one occupied in the conduction state. A conduction electron becomes bound at the defect, ceasing to contribute to the conductivity. This process is very efficient in insulators, so the few conduction electrons provided by impurities and thermal fluctuations are usually trapped at other defects. By definition, an insulator is a solid that does not provide a stable environment for conduction electrons.

Metals have a high density of conduction electrons. The aluminum atom has three valence electrons in a partially filled outer shell. In metallic aluminum the three valence electrons per atom become conduction electrons. The number of conduction electrons is constant, depending on neither temperature nor impurities. Metals conduct electricity at all temperatures, but for most metals the conductivity is best at low temperatures. Divalent atoms, such

Conduction
electrons

Ion
hopping

as magnesium or calcium, donate both valence electrons to become conduction electrons, while monovalent atoms, such as lithium or gold, donate one. As will be recalled, the number of conduction electrons alone does not determine conductivity; it depends on electron mobility as well. Silver, with only one conduction electron per atom, is a better conductor than aluminum with three, for the higher mobility of silver compensates for its fewer electrons.

In metals such as sodium and aluminum, the atoms donate all their valence electrons to the conduction band. The resulting ions are small, occupying only 10–15 percent of the volume of the crystal. The conduction electrons are free to roam through the remaining space. A simple model, which often describes well the properties of the conduction electrons, treats them as interacting neither with the ions nor with each other. The electrons are approximated as free particles wandering easily through the crystal. This concept was first proposed by the German scientist Arnold Johannes Wilhelm Sommerfeld. It works quite well for those metals, known as simple metals, whose conduction electrons are donated from *sp*-shells—for example, aluminum, magnesium, calcium, zinc, and lead. They are called simple because they are aptly described by the simple theory of Sommerfeld.

The transition metals are found in three rows of the periodic table: the first row consists of scandium through nickel, the second row is yttrium through palladium, and the third row is lanthanum plus hafnium through platinum. Within these rows, as the atomic number increases, the electrons fill *d*-states in the outer shell of the atom. In crystal form the transition metal atoms are metals with interesting properties. The *d*-electrons are more tightly bound to the ion centre than are *sp*-electrons. While the *sp*-valence electrons become conduction electrons that move freely through the crystal, the *d*-electrons tend to stay localized near the ion. Neighbouring ions may covalently bond *d*-electrons. In most cases, these *d*-states are only partially filled. Electrons in these *d*-states can conduct as well as those in the *sp*-states, but the electron motion in the *d*-states is not well approximated by the Sommerfeld model of free particles. Instead, the electrons move from ion to ion through the shared covalent bonds of the *d*-electrons. These metals have some conduction electrons donated from *sp*-states and others from *d*-states; therefore, some electrons move freely according to the Sommerfeld model, while others move through the bonds. Each electron switches back and forth between these two modes of conduction, resulting in electron motion that is quite complicated.

An applied voltage causes the electrons of metals to accelerate and contribute to the electric current. The electrons scatter occasionally from imperfections in the crystal, and the rate of scattering determines the mobility. The electrons do not scatter from the ions in the crystal that are located at the expected site in the crystal lattice. The electrons move to accommodate the host ions rather than scatter from them. If an ion is missing, misplaced, or of a different species, however, the electron will scatter from this defect. Ions vibrate around their lattice site, with the amplitude of vibration increasing with temperature. The vibration may cause the ion to be displaced from its crystal site, providing a defect from which an electron will scatter. The resistivity of metals increases at high temperature, owing to the increase in vibrations of the ions in the crystal and the resulting increase in scattering.

Semiconductors have conducting properties intermediate to those of insulators and metals. In some cases the semiconductors are insulators, while in others they are metals. Semiconductors share with insulators the property that they have no conduction electrons in a perfect crystal without thermal fluctuations. Conduction electrons are provided by electrons from impurities or by thermal fluctuation of electrons from atomic shells. The important difference between insulators and semiconductors is in the nature of the traps. A trap is a local electron energy state at a defect. Although the traps in insulators bind conduction electrons tightly, those in semiconductors only weakly bind the electrons. A trapped conduction electron in a semiconductor can be kicked back to the conduction

band by thermal fluctuations. At room temperature, the majority of extra electrons are found in the conduction band rather than in traps. The inability of traps to keep electrons is the main difference between semiconductors and insulators. A semiconductor at room temperature has a sufficient number of conduction electrons to provide good electrical conductivity. Since the mobility of electrons in many semiconductors is exceptionally high, even a small number of conduction electrons is generally sufficient to allow high conductivity.

Phosphorus has five valence electrons, while silicon has four. When a phosphorus atom substitutes for an atom in a silicon crystal lattice, four of its five valence electrons enter covalent bonds. The fifth one is extra, sitting in a shallow trap around the phosphorus site. It is easily excited, however, to the conduction band by thermal fluctuations. At room temperature, there is nearly one conduction electron in silicon for each phosphorus impurity. By controlling the number of impurities, it is possible to control the conductivity of silicon. Other substitutional atoms such as arsenic and antimony also serve as electron donors to the conduction band of silicon.

If a sufficient number of conduction electrons are added to a semiconductor through the introduction of impurities, the electrical properties become metallic. There is a critical concentration of impurities N_c , which depends on the type of impurity. For impurity concentrations less than the critical amount N_c , the conduction electrons become bound in traps at extremely low temperatures, and the semiconductor becomes an insulator. For a concentration of impurities higher than N_c , the conduction electrons are not bound in traps at low temperatures, and the semiconductor exhibits metallic conduction. For phosphorus impurities in silicon, $N_c = 2 \times 10^{18}$ impurities per cubic centimetre. Although this number seems large, it represents about one phosphorus atom for each 100,000 silicon atoms. On a percentage basis, a small number of phosphorus atoms will change silicon from an insulator to a metallic conductor. Other semiconductors have similar properties. In gallium arsenide the critical concentration of impurities for metallic conduction is 100 times smaller than in silicon.

Gallium atoms, like those of phosphorus, can be used as substitutional impurities in silicon. Each atom contributes three electrons to covalent bonds. Since four electrons are needed to complete a tetrahedral arrangement, there is one electron absent per gallium atom from a full set of covalent bonds. The missing electron is called a hole. Holes can move around the crystal in a process similar to the motion of ion vacancies, except in this case there is an electron vacancy. An electron from a nearby covalent bond can jump over and fill the empty electron state, thereby moving the hole to the neighbouring bond. The hole contributes a positive charge, since it is the absence of an electron. The mobility of holes in response to an external voltage is almost as high as the mobility of conduction electrons. A semiconductor may have a high density of impurities that cause holes, and a high electrical conductivity is created by their motion. A *p*-type semiconductor is one with a preponderance of holes; an *n*-type semiconductor has a preponderance of conduction electrons. The symbols *p* and *n* come from the sign of the charge of the particles: *positive* for holes and *negative* for electrons.

Thermal fluctuations can excite an electron out of a covalent bond, making it a conduction electron. The bond is left with a missing electron, which constitutes a hole. Thermal fluctuations thus make electron-hole pairs. Usually the electron and hole separate in space, and each wanders away. The Swiss-American scientist Gregory Hugh Wannier first suggested that the electron and hole could bind together weakly. This bound state, called a Wannier exciton, does exist; the hole has a positive charge, the electron has a negative charge, and the opposites attract. The exciton is observed easily in experiments with electromagnetic radiation. It lives for only a short time—between a nanosecond and a microsecond—depending on the semiconductor. The short lifetime is due to the preference for the electron to reenter a covalent bond state, thereby eliminating both the hole and the conduction electron. This

Sommerfeld model of metals

Conducting properties of semiconductors

Conduction by holes

recombination of electron and hole is easily accomplished from the exciton state, since the two particles are spatially nearby. If the electron and hole escape the exciton state by thermal fluctuation, they travel away from each other. Recombination is then less probable, since it occurs only when the wandering particles pass close to one another again. Recombination also can occur at defect sites. First, one particle becomes bound to the defect, followed by the second particle. The electron and hole are again close to one another, and the electron can reoccupy the covalent bond.

As in metals, the mobility of electrons in semiconductors is limited by electron scattering. For crystals with few defects, the mobility is limited by defect scattering at the lowest temperatures and by ion vibrations at moderate and high temperatures. Since semiconductors with few defects have a small number of conduction electrons, the resistivity is high. The number of conduction electrons is increased in semiconductors by adding impurities. Unfortunately, this also increases the scattering from impurities, which reduces the mobility. Figure 14 shows the resistivity of silicon at room temperature ($T = 300\text{ K}$) as a function of the concentration of impurities. The two curves represent conduction by electrons and by holes. Each grid mark on the graph is a factor of 10. The resistivity varies by a factor of one million from the lowest to the highest concentration of impurities.

Semiconductors with few impurities are good photoconductors. Photoconductivity is the phenomenon in which the electrical conductivity of a solid is increased by exposing it to light. Light is electromagnetic radiation within a specific narrow band of frequencies. The quanta of light are absorbed by the semiconductor, creating electron-hole pairs that provide the electrical conduction. More intense light produces more electron-hole pairs and gives rise to better conductivity. Each semiconductor absorbs light over a specific frequency range, so different semiconductors are used as photoconductors for different ranges of frequency.

Zinc oxide (ZnO) is an interesting material with respect to conductivity. It crystallizes in the wurtzite structure, and its bonding is a mix of ionic and covalent. High-purity single crystals are insulators. Zinc oxide is the most piezoelectric of all materials and is widely used as a transducer in electronic devices. (Piezoelectricity is the property of a crystal to become polarized when subjected to pressure.) Zinc oxide is a good semiconductor when aluminum impurities are included in the crystal. Polycrystalline ceramics of semiconducting zinc oxide conduct well and obey Ohm's law. The addition of small amounts of other oxides, such as those of barium and chromium, causes zinc oxide ceramics to have very nonohmic electrical properties; the electrical current in such ceramics is the most nonlinear of any known material. The current I becomes propor-

From H.H. Landolt and R. Bornstein, *Zahlenwerte und Funktionen aus Naturwissenschaften und Technik* (Numerical Data and Functional Relationships in Science and Technology), new series group III *Kristall- und Festkörperphysik* (Crystal and Solid State Physics), vol. 17a, editor K.H. Hellwege, © Springer-Verlag, Berlin, 1982

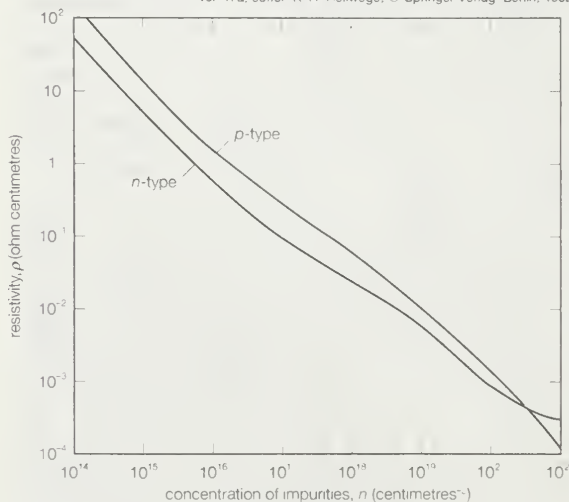


Figure 14: The electrical resistivity (ρ) of a silicon semiconductor at a temperature of 300 K as a function of the concentration of impurities (n).

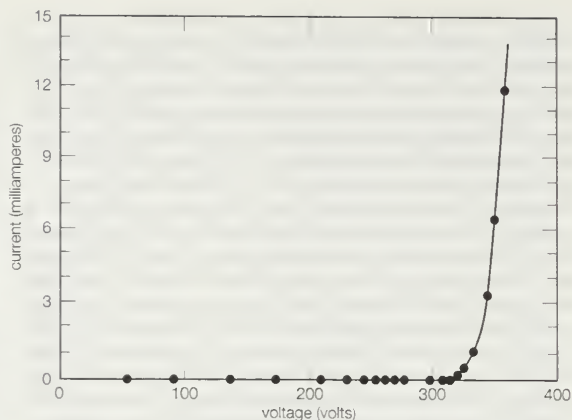


Figure 15: The dependence of current on voltage in a varistor of zinc oxide. The nonlinear behaviour is evident in the sharp upturn in current when the voltage exceeds 330 volts.

W.N. Schultz, General Electric Company

tional to a power of the voltage V^n , where the exponent n has values of more than 100 in certain ranges of voltage. This material is called a varistor, which is a contraction of the words variable and resistor. Zinc oxide varistors are widely used as circuit elements to protect against voltage surges. Figure 15 shows a graph of current versus voltage for a zinc oxide varistor used in household electronics. There is little current until a critical voltage of about 330 volts is reached, at which point the current rises steeply in a nonlinear fashion. Another interesting application of zinc oxide was its former use as a white pigment in paint. It has been replaced by titanium dioxide (TiO_2), however, which is whiter.

MAGNETISM

Electrons are perpetually rotating, and, since the electron has a charge, its spin produces a small magnetic moment. Magnetic moments are small magnets with north and south poles. The direction of the moment is from the south to the north pole. In nonmagnetic materials the electron moments cancel, since there is random ordering to the direction of the electron spins. Whenever two electrons have their moments aligned in opposite directions, their effects tend to cancel. Magnets are formed when a large number of the electrons align their individual moments in the same direction. Only a small percentage of crystals are magnetic. The forces that tend to align the electron spins are subtle. There are three separate parts to the explanation of magnetism. They are as follows:

Formation of magnets

1. Most magnets are composed of atoms whose valence electrons are in d - or f -shells. Atomic shell notation refers to angular momentum, where s has zero unit, p has one, d has two, and f has three. Electrons in d -shells tend to be bound to the ion, and those in f -shells are bound even more tightly.

2. Each electron orbital can be occupied by two electrons—one with spin up and one with spin down. The d -shell has five orbital states and 10 electrons when filled; the f -shell has seven orbital states and 14 electrons when filled. Electrons are added one at a time to the d -states according to the empirical rule that the electrons arrange themselves in the state with the maximum spin and the maximum magnetic moment. If the first electron has spin up, the next four will also have spin up. A maximum of five electrons with spin up are allowed in the d -shell, so the sixth must have spin down. Similarly, the f -shell accepts seven electrons with the same direction of spin before taking electrons with the opposite spin orientation. The order in which electrons fill atomic shells is described by Hund's rules, of which the first is maximizing the total spin. Atoms with electrons in partially filled d - and f -shells usually have a nonzero total spin and thus a net magnetic moment. These magnetic ions are the building blocks for magnetic crystals.

Hund's first rule is due to a phenomenon called electron exchange. As discussed above, a fundamental rule of quantum mechanics, the Pauli exclusion principle, states

Photo-conductivity

that no two electrons with the same direction of spin can occupy the same point in space at the same time. Electrons have charge and repel one another. If two electrons come close together, a large amount of repulsive energy is produced. Physical systems prefer the state of lowest energy, and so electrons avoid such close approach. When their spins are parallel, electrons avoid each other because of the Pauli principle. Electrons in the same shell thus prefer to have their spins parallel, since this configuration keeps the electrons apart and thereby reduces the amount of repulsive energy. The concept of electron exchange is the basis of magnetism. It explains why ions such as iron have large magnetic moments. In divalent iron (Fe^{2+}) the six d -electrons are arranged to achieve maximum electron spin and magnetic moment.

3. Individual ions with fixed magnetic moments may cooperatively align their moments, resulting in the presence of magnetic properties of the crystal as a whole. Ferromagnet crystals have the magnetic moments from all their constituent ions aligned in the same direction; the magnetic moment of the crystal is the summation of the individual moments of the ions. There must be a magnetic force between the different ions that causes them to cooperatively align their moments. This force is also due to electron exchange. The d -orbitals from neighbouring ions overlap weakly into covalent bonds. The d -electrons on the separate ions are shared with the neighbour through covalent bonding. The electron exchange will tend to align the spins on the two neighbours. Aligning all pairs of neighbours aligns all ions. The exchange force between neighbours is much weaker than the force within the atomic shell of one ion. Although weak, the force is sufficient to cause ferromagnetism.

Iron is a typical ferromagnet. Not all bars of iron are magnets; the existence of magnetism is determined by the nature of the domains within the bar. A domain is a region of a crystal in which all the ions are ferromagnetically aligned in the same direction. A bar may be composed of many domains, each having a different magnetic orientation. Such a bar would not appear to be magnetic. Each piece of the bar is magnetic, but the domains have moments that point in different directions, so the bar has no net moment. If the bar of iron is placed in a strong magnetic field, however, the bar becomes magnetic. The field causes the bar to become a single domain with all moments aligned along the external field. The domains do not rotate their moments; instead, the walls between domains move. The domain with a moment along the field grows, while the others become smaller. If removed from the magnetic field, the iron bar will remain magnetized for a considerable time period. Nearly all bars of iron are polycrystalline: they have many small grains of single crystals, which are packed together with random orientation. A grain could be a single domain, a domain could include many grains, or a large grain could have several domains.

Ferromagnetic materials change their magnetic ordering at a characteristic temperature T_c called the Curie temperature. The Curie temperatures for three common ferromagnetics—iron, cobalt, and nickel—are 1,043 K; 1,394 K; and 631 K, respectively. For temperatures below T_c the magnetic moments of the ions are aligned and the crystal is magnetic. For temperatures above T_c the crystal is not ferromagnetic, since the individual atomic moments are no longer aligned. Above T_c the moments have short-range order but not long-range order. Short-range order means there is local ordering. If a moment points in one direction, its neighbours have a tendency to point in the same direction. This tendency is maintained over several lattice sites but is not maintained for long distances. Long-range order is the tendency for moments to align for large distances. For temperatures a few degrees below T_c the moments have strong short-range order but only a small amount of long-range order, so the bar is not very magnetic. The tendency for long-range order increases at lower temperature. The Curie temperature is the point where long-range order begins as the temperature is lowered. The magnetization of nickel as a function of temperature is shown in Figure 16.

If an iron bar is heated to a temperature above T_c the

Ferro-
magnetic
materials

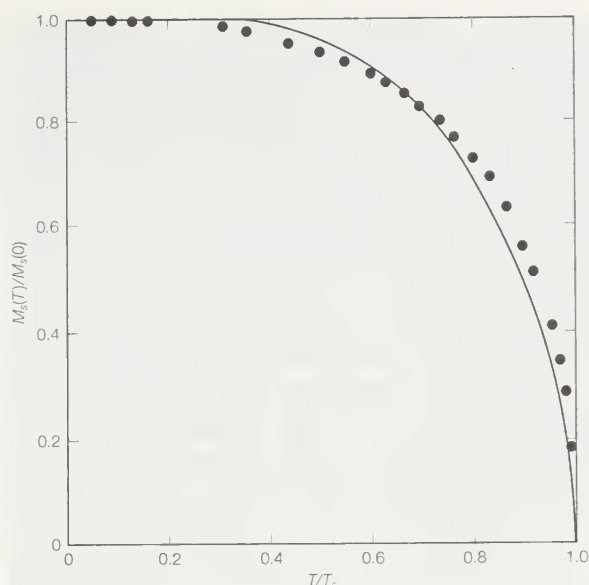


Figure 16: The total magnetization in nickel as a function of temperature. There is no long-range order above the Curie temperature T_c . $M_s(T)$ is the magnetization at temperature T , and $M_s(0)$ is the magnetization at absolute zero.

From P. Weiss and R. Forrer, *Annalen der Physik*, 5,153 (1926) in C. Kittel, *Introduction to Solid State Physics*, 4th ed., copyright © 1971 John Wiley & Sons, Inc.; reprinted by permission of John Wiley & Sons, Inc.

bar is no longer magnetic. If the bar is then cooled to a temperature below T_c the grains become magnetic, but they orient their moments in random directions, so the bar as a whole is not magnetic. A bar can be demagnetized by heating the bar and then cooling it. By inserting it in a large magnetic field, the bar can be remagnetized.

Ferromagnetism is found in many insulators as well as metals. Chromium bromide (CrBr_3) is an insulator since chromium is trivalent and a bromine atom needs one electron to complete its outer shell. The trivalent chromium atoms each have a moment, and these align ferromagnetically below the Curie temperature of 37 K. Gadolinium chloride (GdCl_3 ; $T_c = 2.2$ K) and europium oxide (EuO ; $T_c = 77$ K) are two other examples among many.

Many crystals have magnetic ions that are ordered in arrangements other than ferromagnetic. In antiferromagnetic ordering, the moments pointing in one direction are balanced by others pointing in the opposite direction, with the result that the substance has no net magnetization. The exchange interaction between ions in this case has the opposite sign and favours the alternate arrangements of spins. The sign of the exchange interaction between ions depends on the length of the covalent bond and the bonding angles; it may have either orientation. The characteristic temperature associated with antiferromagnetism is called the Néel temperature T_N . Below T_N the ions are antiferromagnetically ordered, while above this temperature there is no long-range antiparallel order. Some examples of antiferromagnetic crystals are manganese oxide (MnO ; $T_N = 116$ K), manganese sulfide (MnS ; $T_N = 160$ K), and iron oxide (FeO ; $T_N = 198$ K). Manganese oxide is an insulator since manganese atoms are divalent and oxygen atoms accept two electrons. The manganese ion has a fixed magnetic moment. The crystal structure of manganese oxide is the same as that of sodium chloride shown in Figure 7B. Below the Néel temperature the atomic unit cell doubles in size to include two atoms of each type of ion. This is necessary because below T_N neighbouring manganese atoms have moments in the opposite direction and are no longer equivalent; the unit cell must therefore include one moment in each of the two directions. Fluorides such as manganese fluoride (MnF_2), iron (II) fluoride (FeF_2), cobalt fluoride (CoF_2), and nickel fluoride (NiF_2) are other crystals that exhibit antiferromagnetic ordering of the transition metal ions.

Antiferro-
magnetic
materials

Ferrimagnetism is another type of magnetic ordering. In ferrimagnets the moments are in an antiparallel alignment, but they do not cancel. The best example of a ferrimagnetic mineral is magnetite (Fe_3O_4). Two iron ions

are trivalent, while one is divalent. The two trivalent ions align with opposite moments and cancel one another, so the net moment arises from the divalent iron ion. The historic lodestone that was the first magnetic material discovered was a form of magnetite. Another class of ferromagnets has the garnet structure; $Y_3Fe_5O_{12}$ is one such crystal. Only the iron ions are magnetic. Three point in one direction and two in the other, so there is a net magnetic moment of one iron ion in the unit cell. However, if a rare-earth ion such as gadolinium (Gd) is substituted for yttrium (Y), then the rare-earth ion also contributes to the ferromagnetism.

Ferrites

Ferrites are oxides of iron with the formula MFe_2O_4 , where M is a divalent ion such as nickel, zinc, cadmium, manganese, or magnesium. The ferric iron ion is trivalent, and the oxygen ion accepts two electrons. Actually M can also be divalent iron, forming magnetite (Fe_3O_4). The crystal structure is called spinel, which is the mineral name for $MgAl_2O_4$. Ferrites are electrical insulators with magnetic ordering. Their insulating quality makes them useful as magnetic cores. When metallic ferromagnetic materials are exposed to alternating magnetic fields, significant heating losses occur from eddy currents. Ferrite magnets greatly reduce such heat losses because of their high resistivity. They also absorb electromagnetic radiation of very long wavelengths. This property is unusual, since metals reflect such radiation while insulators transmit it. In small crystallites of ferrites, the electromagnetic radiation causes the magnetic moment to rotate at the frequency of light. The absorption frequency of the light depends on the detailed shape of the crystal grain. A polycrystalline material, with a range of grain sizes and shapes, absorbs over a broad range of radar frequencies. Ferrite crystals are a major ingredient in snoop paint, which makes stealth airplanes undetectable by radar. Regular airplanes reflect radar, and the reflected signals can be used to locate the aircraft. Stealth planes absorb the radar signals and so cannot be located in this way.

Magnetic ions have interesting properties when they are found as impurities in nonmagnetic crystals. They usually retain their magnetic moment, so small magnets are distributed randomly throughout the crystal. If the host crystal is a metal, the magnetic impurities make an interesting

contribution to the electrical resistivity. The conduction electrons scatter from the magnetic impurity. Since the conduction electron and the impurity both have spin, they can mutually flip spins while scattering. The spin-flip scattering is strong at low temperatures and actually increases slightly as temperature decreases. This phenomenon is called the Kondo effect after the Japanese theoretical physicist Jun Kondo, who first explained the increase in resistivity resulting from magnetic impurities. There is a characteristic temperature, called the Kondo temperature, which depends on the impurity and on the metallic host. The resistivity increases at low temperature, starting near the Kondo temperature. A typical example of a Kondo system is iron impurities in copper: the system's Kondo temperature is 24 K. The solid line in Figure 17 shows the resistivity in copper at low temperature when there are 110 iron impurities per 1,000,000 copper atoms. The dashed line *a* is the resistivity in the absence of impurities. It increases at higher temperature because the electron scatters from ion vibrations. The dashed line *b* is the resistivity from spin-flip scattering. Nearly all transition metal atoms are found as magnetic impurities in copper or gold. Each system has a different Kondo temperature, which varies from 1,000 K to a fraction of 1 K. The spin-flip part of the electrical resistivity is unique in that it is large at low temperatures and decreases at high temperatures; most contributions to the resistivity increase at high temperatures. (G.D.Ma.)

Quasicrystals

Quasicrystals are metal alloys whose novel symmetries challenge the traditional dogma of crystallography. Although they appear at casual inspection to be ordinary crystals, they are not. Their symmetries elude the classification scheme of crystal structures, which enumerates the combinations of translational and rotational symmetries that are allowed according to the laws of crystallography. While ordinary crystals place atoms in periodic lattices, quasicrystals arrange atoms in a quasiperiodic fashion. Although these structures surprised the scientific community, it now appears that quasicrystals rank among the most common structures in alloys of aluminum with such metals as iron, cobalt, or nickel. While no major commercial applications yet exploit properties of the quasicrystalline state directly, quasicrystals form in compounds noted for their high strength and light weight, suggesting potential applications in aerospace and other industries.

STRUCTURE AND SYMMETRY

Dan Shechtman, a researcher from Technion, a part of the Israel Institute of Technology, and his colleagues at the National Bureau of Standards (now the National Institute of Standards and Technology) in Gaithersburg, Md., discovered quasicrystals in 1984. A research program of the U.S. Air Force sponsored their investigation of the metallurgical properties of aluminum-iron and aluminum-manganese alloys. Shechtman and his coworkers mixed aluminum and manganese in a roughly six-to-one proportion and heated the mixture until it melted. The mixture was then rapidly cooled back into the solid state by dropping the liquid onto a cold spinning wheel, a process known as melt spinning. When the solidified alloy was examined using an electron microscope, a novel structure was revealed. It exhibited fivefold symmetry, which is forbidden in crystals, and long-range order, which is lacking in amorphous solids. Its order, therefore, was neither amorphous nor crystalline. Many other alloys with these same features have subsequently been produced.

The electron microscope has played a significant role in the investigation of quasicrystals. It is a versatile tool that can probe many important aspects of the structure of matter. Low-resolution scanning electron microscopy magnifies the shapes of individual grains. Grains of a quasicrystalline aluminum-copper-iron alloy imaged with this technique are shown in Figure 18. Symmetries of solid grains often reflect the internal symmetries of the underlying atomic positions. Grains of salt, for example, take cubical shapes consistent with the cubic symmetries

Discovery
of quasi-
crystals

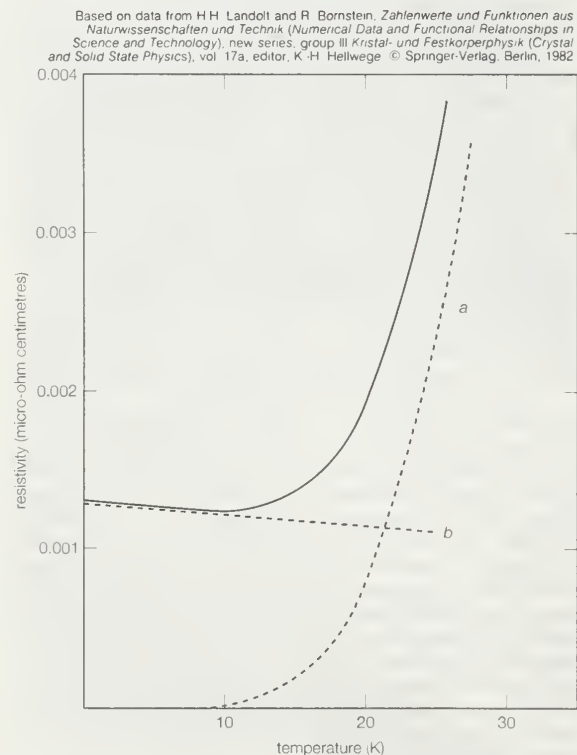


Figure 17: The electrical resistivity (solid line) of copper at low temperature when there are 110 iron atoms per 1,000,000 copper atoms. Without impurities the resistivity curve would be *a*. Curve *b* is the resistivity from the spin-flip scattering that results from the iron impurities.

of their crystal lattices. The shape observed in Figure 18 is called a pentagonal dodecahedron. Its 12 faces are regular pentagons, with axes of fivefold rotational symmetry passing through them. That is to say, rotations about this axis by 72° leave the appearance of the grain unchanged. In a full 360° rotation the grain will repeat itself in appearance five times, once every 72° . There are also axes of twofold rotational symmetry passing through the edges and axes of threefold rotational symmetry passing through the vertices. This is also known as icosahedral symmetry because the icosahedron is the geometric dual of the pentagonal dodecahedron. At the centre of each face on an icosahedron, the dodecahedron places a vertex, and vice versa. The symmetry of a pentagonal dodecahedron or icosahedron is not among the symmetries of any crystal structure, yet this is the symmetry that was revealed in the electron microscope image of the aluminum-manganese alloy produced by Shechtman and his colleagues.

High-resolution electron microscopy magnifies to such a great degree that patterns of atomic positions may be determined. In ordinary crystals such a lattice image reveals regularly spaced rows of atoms. Regular spacing implies spatial periodicity in the placement of atoms. The angles between rows indicate rotational symmetries of the atomic positions. A high-resolution electron microscope image of quasicrystalline aluminum-manganese-silicon is shown in Figure 19. Rows of atoms may be visualized by glancing along the page. Parallel rows occur in five sets, rotated from one another by 72° , confirming that the fivefold symmetry suggested by the shape of the pentagonal dodecahedron grain reflects a fivefold symmetry in the actual placement of atoms.

F.W. Gayle, *Journal of Metals*, vol. 40, no. 5, May 1988

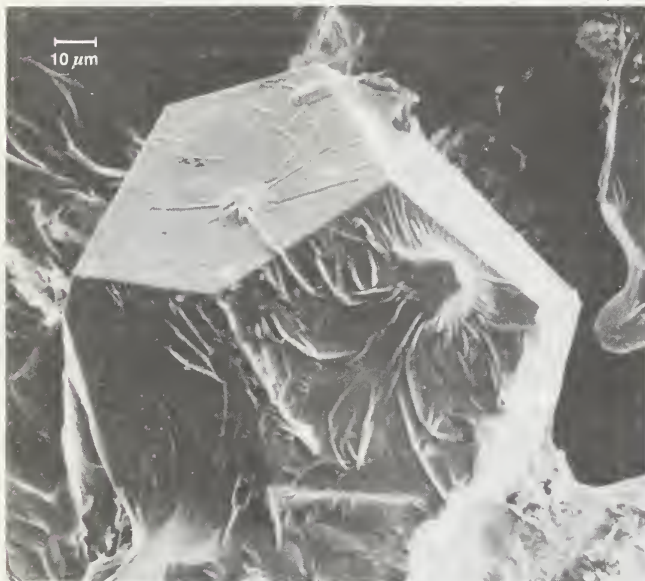


Figure 18: A scanning electron microscope image of quasicrystalline aluminum-copper-iron, revealing the pentagonal dodecahedral shape of the grain.

Fivefold symmetry axes are forbidden in ordinary crystals, while other axes, such as sixfold axes, are allowed. The reason is that translational periodicity, which is characteristic of crystal lattices, cannot be present in structures with fivefold symmetry. Figures 20 and 21 can be used to illustrate this concept. The triangular array of atoms in Figure 20 has axes of sixfold rotational symmetry passing through each atomic position. The arrows represent translational symmetries of this crystalline structure. That is, if the entire array of atoms is displaced along one of these arrows, say the one labeled *a*, all new atomic positions coincide with the locations of other atoms prior to the displacement. Such a displacement of atoms that leaves atomic positions invariant is called a symmetry of the crystal. In Figure 20, if two different symmetries are combined such that the structure is first displaced along arrow *a* and then along arrow *b*, the net result is equivalent to a displacement along arrow *c*, which itself must be a symmetry of the structure. Again, atomic sites coincide before

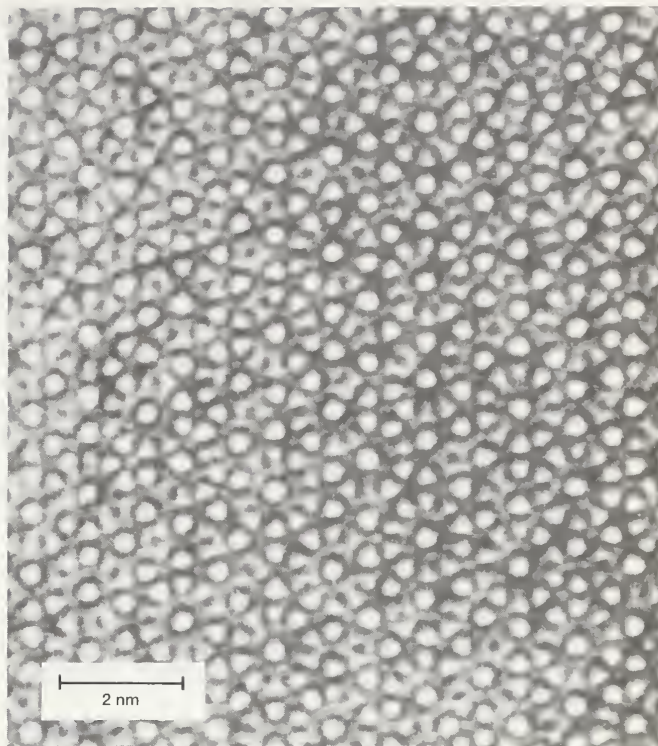


Figure 19: A high-resolution electron microscope image of quasicrystalline aluminum-manganese-silicon, revealing a fivefold symmetry of atomic positions. A glancing view along this figure reveals a Fibonacci sequence of dark and light rows.

Courtesy, Kenji Hiraga

and after the displacement. Repeated displacements along the same arrow demonstrate the translational periodicity of the crystal.

The atomic arrangement shown in Figure 21 exhibits fivefold rotational symmetry but lacks the translational symmetries that must be present in a crystalline structure. The arrows (other than arrow *c*) represent displacements that leave the arrangement invariant. Assume they are the shortest such displacements. Now, as before, consider the combinations of two symmetries *a* and *b* with the net result *c*. The length of *c* is smaller than either *a* or *b* by a factor $\tau = (\sqrt{5} + 1)/2$, which is known as the golden mean. The new atomic position, outlined with a dotted line, does not coincide with a previous atomic position, indicating that the structure does not exhibit translational periodicity. Therefore, an array of atoms may not simultaneously display fivefold rotational symmetry and translational periodicity, for, if it did, there would be no lower limit to the spacing between atoms.

In fact, the compatibility of translational periodicity with sixfold rotational symmetry (as shown in Figure 20) is a remarkable accident, for translational periodicity is not possible with most rotational symmetries. The only

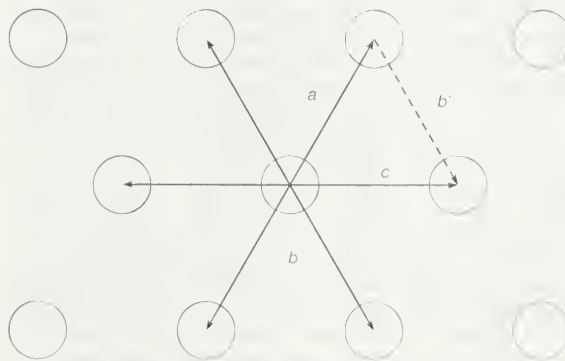


Figure 20: Hexagonal lattice of atomic sites. Arrows indicate translational symmetries of the lattice. Combining two symmetries (*a* and *b*) produces a third (*c*).

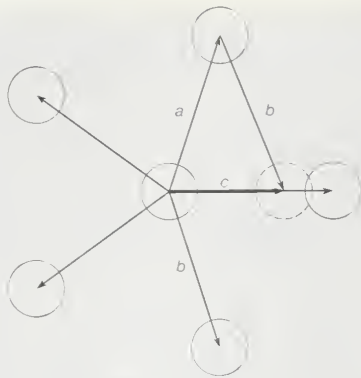


Figure 21: *Pentagonal arrangement of atoms.* Arrows indicate hypothetical shortest translational symmetries. Combining two displacements produces a new displacement that is shorter than the hypothetical symmetries.

Symmetries found in quasicrystals

allowed symmetry axes in periodic crystals are twofold, threefold, fourfold, and sixfold. All others are forbidden owing to the lack of minimum interatomic separation. In particular, fivefold, eightfold, tenfold, and twelvefold axes cannot exist in crystals. These symmetries are mentioned in particular because they have been reported in quasicrystalline alloys.

Since the aluminum-manganese-silicon quasicrystal shown in Figure 19 clearly reveals an axis of fivefold symmetry, it may be concluded that the arrangement of atoms lacks translational periodicity. That, in itself, is no great surprise, for many materials lack translational periodicity. Amorphous metals, for example, are frequently produced by the same melt-spinning process that was employed in the discovery of quasicrystals. Amorphous metals have no discrete rotational symmetries, however, and high-resolution electron microscope images reveal no rows of atoms. The arrangement of atoms in a quasicrystal displays a property called long-range order, which is lacking in amorphous metals. Long-range order permits rows of atoms to span Figure 19 and maintains agreement of row orientations across the figure. Ordinary crystal structures, such as that of Figure 20, display long-range order. Strict rules govern the relative placement of atoms at remote locations in solids with long-range order.

Electron diffraction confirms the presence of long-range order in both crystals and quasicrystals. Quantum mechanics predicts that particles such as electrons move through space as if they were waves, in the same manner that light travels. When light waves strike a diffraction grating, they are diffracted. White light breaks up into a rainbow, while monochromatic light breaks up into discrete sharp spots. Similarly, when electrons strike evenly spaced rows of atoms within a crystalline solid, they break up into a set of bright spots known as Bragg diffraction peaks. Symmetrical arrangements of spots reveal axes of rotational symmetry in the crystal, and spacings between the discrete spots relate inversely to translational periodicities. Amorphous metals contain only diffuse rings in their diffraction patterns since long-range coherence in atomic positions is required to achieve sharp diffraction spots.

The original electron diffraction pattern of quasicrystalline aluminum-manganese published by Shechtman and his coworkers is shown in Figure 22. Rings of 10 bright spots indicate axes of fivefold symmetry, and rings of six bright spots indicate axes of threefold symmetry. The twofold symmetry axes are self-evident. The angles between these axes, indicated on the figure, agree with the geometry of the icosahedron. The very existence of spots at all indicates long-range order in atomic positions. Recalling the earlier result that fivefold symmetry axes are forbidden in crystalline materials, a paradox is presented by quasicrystals. They have long-range order in their atomic positions, but they must lack spatial periodicity.

Dov Levine and Paul Steinhardt, physicists at the University of Pennsylvania, proposed a resolution of this apparent conflict. They suggested that the translational order

of atoms in quasicrystalline alloys might be quasiperiodic rather than periodic. Quasiperiodic patterns share certain characteristics with periodic patterns. In particular, both are deterministic—that is, rules exist that specify the entire pattern. These rules create long-range order. Both periodic and quasiperiodic patterns have diffraction patterns consisting entirely of Bragg peaks. The difference between quasiperiodicity and periodicity is that a quasiperiodic pattern never repeats itself. There are no translational symmetries, and, consequently, there is no minimum spacing between Bragg peaks. Although the peaks are discrete, they fill the diffraction pattern densely.

The most well-known quasiperiodic pattern may be the Fibonacci sequence, discovered during the Middle Ages in the course of studies conducted on rabbit reproduction. Consider the following rules for birth and maturation of rabbits. Start with a single mature rabbit (denoted by the symbol L for large) and a baby rabbit (denoted by S for small). In each generation every L rabbit gives birth to a new S rabbit, while each preexisting S rabbit matures into an L rabbit. A table of rabbit sequences may be established as follows. Start with an L and an S side by side along a line. Replace the L with LS and the S with L to obtain LSL and repeat the procedure as shown in Table 1. The numbers of rabbits present after each generation are the Fibonacci numbers. The population grows exponentially over time, with the population of each generation approaching τ (the golden mean) multiplied by the population of the previous generation. The sequence of L and S symbols forms a quasiperiodic pattern. It has no subunit that repeats itself periodically. In contrast, a periodic sequence such as LSLLSLSLSLSLSL . . . has a fundamental unit (LSL) that is precisely repeated at equal intervals. In crystallography such a repeated unit is called a unit cell. Quasiperiodic sequences have no unit cell of finite size. Any portion of the Fibonacci sequence is repeated in-

Quasiperiodicity

(Left) D. Shechtman, *Physical Review Letters*, vol. 53 no. 20 Nov. 1984

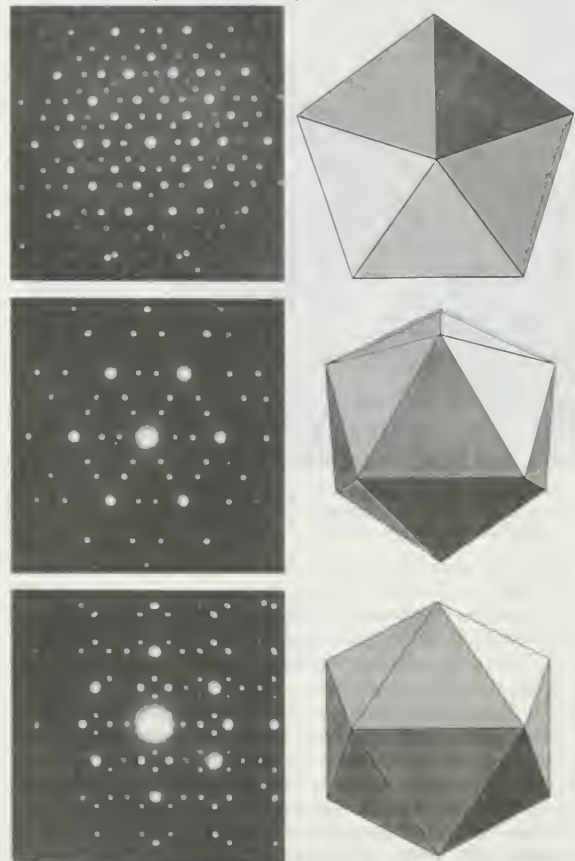


Figure 22: (Left) Electron diffraction patterns of quasicrystalline aluminum-manganese. (Top left) View is along the fivefold symmetry axis; (centre left) rotating by 37.38° reveals the threefold axis, and (bottom left) rotating by 58.29° reveals the twofold axis. (Right) Corresponding views of icosahedrons show that quasicrystalline symmetries match those of the icosahedrons.

Table 1: Fibonacci Sequences of Rabbits

generation	sequence	mature rabbits	babies
1	LS	1	1
2	LSL	2	1
3	LSLLS	3	2
4	LSLLSLSL	5	3
5	LSLLSLSLSLS	8	5
6	LSLLSLSLSLSLSLSL	13	8

finitely often, but at intervals that are not periodic. These intervals themselves form a Fibonacci sequence.

An example of a two-dimensional pattern that combines fivefold rotational symmetry with quasiperiodic translational order is the Penrose pattern, discovered by the English mathematical physicist Roger Penrose and shown in Figure 23. The diffraction pattern of such a sequence closely resembles the fivefold symmetric patterns of Figure 22. The rhombic tiles are arranged in sets of parallel rows; the shaded tiles represent one such set, or family. Five families of parallel rows are present in the figure, with 72° angles between the families, although only one of the five has been shaded. Within a family the spacings between rows are either large (L) or small (S), as labeled in the margin. The ratio of widths of the large rows to the small rows is equal to the golden mean τ , and the quasiperiodic sequence of large and small follows the Fibonacci sequence. An example of the use of Penrose tilings in an architectural application is shown in Figure 24.

Levine's and Steinhardt's proposal that quasicrystals possess quasiperiodic translational order can be examined in terms of the high-resolution electron micrograph in Figure 19. The rows of bright spots are separated by small and large intervals. As in the Penrose pattern, the length of the large interval divided by the length of the small one equals the golden mean, and the sequence of large and small reproduces the Fibonacci sequence. Levine's and Steinhardt's proposal appears consistent with the electron diffraction results. The origin of the name quasicrystals arises from the fact that these materials have quasiperiodic translational order, as opposed to the periodic order of ordinary crystals.

Figures 18, 19, and 22 represent quasicrystals with the symmetry of an icosahedron. Icosahedral quasicrystals occur in many intermetallic compounds, including aluminum-copper-iron, aluminum-manganese-palladium, aluminum-magnesium-zinc, and aluminum-copper-lithium. Other crystallographically forbidden symmetries have been observed as well. These include decagonal symmetry, which exhibits tenfold rotational symmetry within two-dimensional atomic layers but ordinary translational periodicity perpendicular to these layers. Decagonal symmetry has been found in the compounds aluminum-copper-cobalt and aluminum-nickel-cobalt. Structures that are periodic in two dimensions but follow a Fibonacci sequence in the remaining third dimension occur in aluminum-copper-nickel.

All the compounds named thus far contain aluminum. Indeed, it appears that aluminum is unusually prone to quasicrystal formation, but there do exist icosahedral quasicrystals without it. Some, like gallium-magnesium-zinc, simply substitute the chemically similar element gallium for aluminum. Others, like titanium-manganese, appear chemically unrelated to aluminum-based compounds. Furthermore, some quasicrystals such as chromium-nickel-silicon and vanadium-nickel-silicon display octagonal and dodecagonal structures with eightfold or twelvefold symmetry, respectively, within layers and translational periodicity perpendicular to the layers.

The origin of quasicrystalline order remains in question. No proven explanation clarifies why a material favours crystallographically forbidden rotational symmetry and translational quasiperiodicity when at nearby compositions it forms more conventional crystal structures. The American chemist Linus Pauling noted that these related crystalline structures frequently contain icosahedral motifs within their unit cells, which are then repeated periodically. Pauling proposed that quasicrystals are really ordinary crystalline materials caught out of equilibrium by a

type of crystal defect called twinning, in which unit cells are attached at angles defined by these icosahedral motifs. While this may be a reasonable model for rapidly cooled alloys such as Shechtman's original aluminum-manganese, other compounds, such as aluminum-copper-iron, possess quasicrystalline structures in thermodynamic equilibrium. These quasicrystals can be grown slowly and carefully using techniques for growth of high-quality conventional crystals. The more slowly the quasicrystal grows, the more perfect will be its rotational symmetry and quasiperiodicity. Measuring the sharpness of diffraction pattern spots shows perfect ordering on length scales of at least 30,000 angstroms in these carefully prepared quasicrystals. Twinning cannot account for such long-range order.

Levine and Steinhardt proposed that matching rules, such as those Penrose discovered to determine proper placement of his tiles to fill the plane quasiperiodically, may force the atoms into predefined, low-energy locations. Such a mechanism cannot be the complete explanation, though, since the compound forms ordinary crystalline structures at nearby compositions and temperatures. Indeed, it appears that, when quasicrystals are thermodynamically stable phases, it is only over a limited range of temperatures close to the melting point. At lower temperatures they transform into ordinary crystal structures. Thermodynamics predicts that the stable structure is the one that minimizes the free energy, defined as the ordinary energy minus the product of the temperature and the entropy. It is likely that entropy (a measure of fluctuations around an ideal structure) must be considered in addition to energy to explain stability of quasicrystals.

PROPERTIES

Along with their novel structures and symmetries, quasicrystals are expected to exhibit unusual properties. Both their elastic and their electronic behaviour distinguish quasicrystals from ordinary crystalline metals. Elastic response may be studied by measuring the speed of sound waves propagating through the metal. Sound speeds usually vary depending on the direction of propagation relative to axes of high rotational symmetry. Because the icosahedron has such high symmetry—it is closer to a sphere than is, for instance, a cube—the sound speeds turn out to be independent of the direction of propagation. Longitudinal sound waves (with displacements parallel to the direction of propagation) have speeds different from transverse waves (with displacements perpendicular to the direction of propagation), as is the case for all matter. Because the sound speeds do not depend on direction of propagation, only two elastic constants are required to specify acoustic properties of icosahedral quasicrystals. In contrast, cubic crystals require three elastic constants, and lower-symmetry crystals require up to 21 constants.

As a consequence of the translational quasiperiodicity, there exists a second type of elastic deformation beyond the ordinary sound wave, or phonon. Known as phasons,

Elastic properties

The Penrose pattern

Symmetries observed in quasicrystals

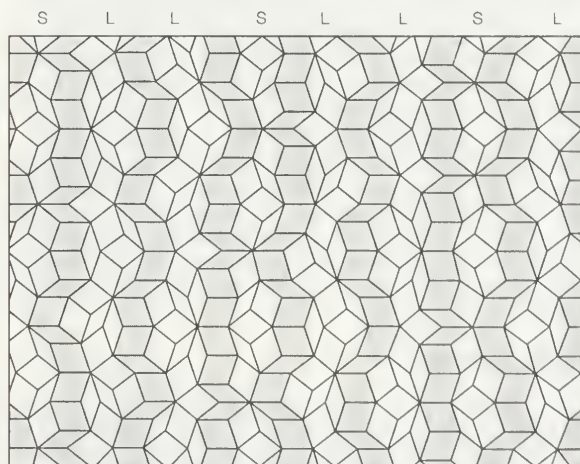


Figure 23: A Penrose tiling.

The plane is covered by rhombuses (deformed squares). Tiles with parallel edges lie in rows (shaded) separated by large (L) and small (S) intervals.

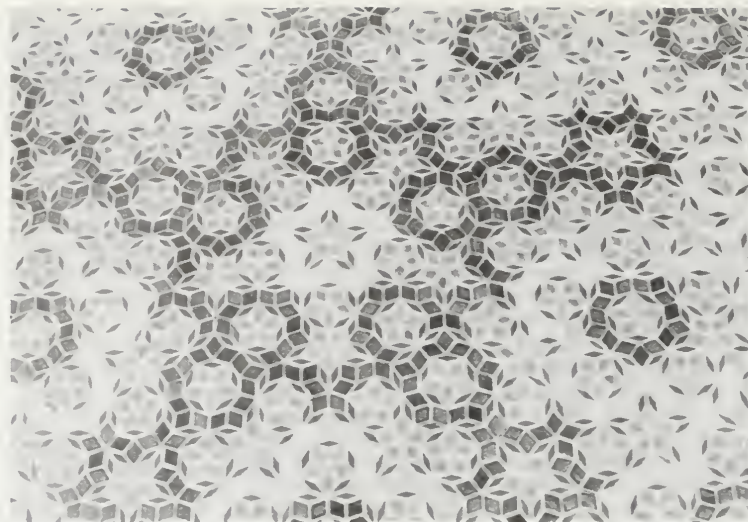


Figure 24: *Penrose tiles.* Penrose tilings are used in architectural applications, as in these ceramic tiles coloured to emphasize the five-pointed stars and bands.

By courtesy of Saxe-Patterson Ceramics for Architecture

these elastic deformations correspond to rearrangements of the relative atomic positions. Removal of a phason requires adjusting positions of all atoms within a row of atoms in a quasicrystalline structure. At low temperatures motion of atoms within the solid is difficult, and phason strain may be easily frozen into the quasicrystal, limiting its perfection. At high temperatures, close to the melting point, phasons continually fluctuate, and atoms jump from place to place.

The electric properties of quasicrystals have proved to be rather unusual. Unlike their constituent elements, which tend to be good electrical conductors, quasicrystals conduct electricity poorly. For alloys of aluminum-copper-ruthenium these conductivities differ by as much as a factor of 100. As the perfection of the quasicrystalline order grows, the conductivity drops. Such behaviour is consistent with the appearance of a gap in the electronic density of states at the Fermi surface, which is the energy level separating filled electronic states from empty ones. Since it is only Fermi-surface electrons that carry current, a vanishingly small density of such electronic states leads to low electrical conductivities in semiconductors and insulators. Such a gap in the density of states may also play a role in explaining the formation of quasicrystalline structures. This is known as the Hume-Rothery rule for alloy formation. Since the Fermi-surface electrons are the highest-energy electrons, diminishing the number of such electrons may lower the overall energy.

The mechanical properties of quasicrystals are especially significant because the desire to develop a material that exhibited these properties motivated the investigators who discovered quasicrystals. Mechanical properties also relate to their first potential practical applications. Quasicrystals are exceptionally brittle. They have few dislocations, and those present have low mobility. Since metals bend by creating and moving dislocations, the near absence of dislocation motion causes brittleness. On the positive side, the difficulty of moving dislocations makes quasicrystals extremely hard. They strongly resist deformation. This makes them excellent candidates for high-strength surface coatings. Indeed, the first successful application of quasicrystals was as a surface treatment for aluminum frying pans.

Liquid crystals

Liquid crystals, their very name an oxymoron, blend structures and properties of the normally disparate liquid and crystalline solid states. Liquids can flow, for example, while solids cannot, and crystalline solids possess special symmetry properties that liquids lack. Ordinary solids melt into ordinary liquids as the temperature increases—*e.g.*, ice melts into liquid water. Some solids actually melt twice

or more as temperature rises. Between the crystalline solid at low temperatures and the ordinary liquid state at high temperatures lies an intermediate state, the liquid crystal. Liquid crystals share with liquids the ability to flow but also display symmetries inherited from crystalline solids. The resulting combination of liquid and solid properties allows important applications of liquid crystals in the displays of such devices as wristwatches, calculators, portable computers, and flat-screen televisions.

STRUCTURE AND SYMMETRY

Crystals exhibit special symmetries when they slide in certain directions or rotate through certain angles. These symmetries can be compared to those encountered when walking in a straight line through empty space. Regardless of the direction or distance of each step, the view remains the same, as there are no landmarks by which to measure one's progress. This is called continuous translational symmetry because all positions look identical. Figure 25A illustrates a crystal in two dimensions. Such a crystal lat-

Electric properties

Translational symmetry

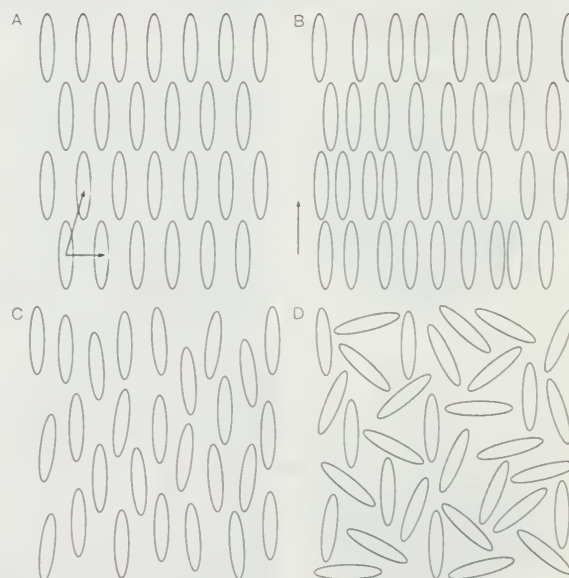


Figure 25: *Arrangements of molecules.* (A) Ordinary crystals break the continuous rotational and translational symmetry of free space; discrete symmetries remain. (B) Smectic liquid crystals show broken translational symmetry in only one direction. (C) Nematic liquid crystals break only rotational symmetry. (D) Isotropic liquids share the continuous translational and rotational symmetry of free space. This symmetry may not be apparent in a single snapshot of molecular positions and orientations

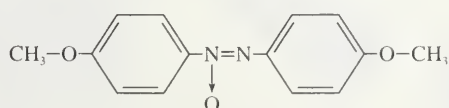
tic breaks the continuous translational symmetry of free space; starting at one molecule there is a finite distance to travel before reaching the next. Some translational symmetry is present, however, because, by moving the proper distance in the proper direction, one is guaranteed to locate additional molecules on repeated excursions. This property is called discrete translational periodicity. The two-dimensional picture of a crystal displays translational periodicity in two independent directions. Real, three-dimensional crystals display translational periodicity in three independent directions.

Rotational symmetries can be considered in a similar fashion. From one point in empty space, the view is the same regardless of which direction one looks. There is continuous rotational symmetry—namely, the symmetry of a perfect sphere. In the crystal shown in Figure 25A, however, the distance to the nearest molecule from any given molecule depends on the direction taken. Furthermore, the molecules themselves may have shapes that are less symmetric than a sphere. A crystal possesses a certain discrete set of angles of rotation that leave the appearance unchanged. The continuous rotational symmetry of empty space is broken, and only a discrete symmetry exists. Broken rotational symmetry influences many important properties of crystals. Their resistance to compression, for example, may vary according to the direction along which one squeezes the crystal. Transparent crystals, such as quartz, may exhibit an optical property known as birefringence. When a light ray passes through a birefringent crystal, it is bent, or refracted, at an angle depending on the direction of the light and also its polarization, so that the single ray is broken up into two polarized rays. This is why one sees a double image when looking through such crystals.

In a liquid such as the one shown in Figure 25D, all the molecules sit in random positions with random orientations. This does not mean that there is less symmetry than in the crystal, however. All positions are actually equivalent to one another, and likewise all orientations are equivalent, because in a liquid the molecules are in constant motion. At one instant the molecules in the liquid may occupy the positions and orientations shown in Figure 25D, but a moment later the molecules will move to previously empty points in space. Likewise, at one instant a molecule points in one direction, and the next instant it points in another. Liquids share the homogeneity and isotropy of empty space; they have continuous translational and rotational symmetries. No form of matter has greater symmetry.

As a general rule, molecules solidify into crystal lattices with low symmetry at low temperatures. Both translational and rotational symmetries are discrete. At high temperatures, after melting, liquids have high symmetry. Translational and rotational symmetries are continuous. High temperatures provide molecules with the energy needed for motion. The mobility disorders the crystal and raises its symmetry. Low temperatures limit motion and the possible molecular arrangements. As a result, molecules remain relatively immobile in low-energy, low-symmetry configurations.

Liquid crystals, sometimes called mesophases, occupy the middle ground between crystalline solids and ordinary liquids with regard to symmetry, energy, and properties. Not all molecules have liquid crystal phases. Water molecules, for example, melt directly from solid crystalline ice into liquid water. The most widely studied liquid-crystal-forming molecules are elongated, rodlike molecules, rather like grains of rice in shape (but far smaller in size). A popular example is the molecule *p*-azoxyanisole (PAA):



Typical liquid crystal structures include the smectic shown in Figure 25B and the nematic in Figure 25C (this nomenclature, invented in the 1920s by the French scientist Georges Friedel, will be explained below). The

smectic phase differs from the solid phase in that translational symmetry is discrete in one direction—the vertical in Figure 25B—and continuous in the remaining two. The continuous translational symmetry is horizontal in the figure, because molecule positions are disordered and mobile in this direction. The remaining direction with continuous translational symmetry is not visible, because this figure is only two-dimensional. To envision its three-dimensional structure, imagine the figure extending out of the page.

In the nematic phase all translational symmetries are continuous. The molecule positions are disordered in all directions. Their orientations are all alike, however, so that the rotational symmetry remains discrete. The orientation of the long axis of a nematic molecule is called its director. In Figure 25C the nematic director is vertical.

It was noted above that, as temperature decreases, matter tends to evolve from highly disordered states with continuous symmetries toward ordered states with discrete symmetries. This can occur through a sequence of symmetry-breaking phase transitions. As a substance in the liquid state is reduced in temperature, rotational symmetry breaking creates the nematic liquid crystal state in which molecules are aligned along a common axis. Their directors are all nearly parallel. At lower temperatures continuous translational symmetries break into discrete symmetries. There are three independent directions for translational symmetry. When continuous translational symmetry is broken along only one direction, the smectic liquid crystal is obtained. At temperatures sufficiently low to break continuous translational symmetry in all directions, the ordinary crystal is formed.

The mechanism by which liquid crystalline order is favoured can be illustrated through an analogy between molecules and grains of rice. Collisions of molecules require energy, so the greater the energy, the greater the tolerance for collisions. If rice grains are poured into a pan, they fall at random positions and orientations and tend to jam up against their neighbours. This is similar to the liquid state illustrated in Figure 25D. After the pan is shaken to allow the rice grains to readjust their positions, the neighbouring grains tend to line up. The alignment is not perfect across the sample owing to defects, which also can occur in nematic liquid crystals. When all grains align, they have greater freedom to move before hitting a neighbour than they have when they are disordered. This produces the nematic phase, illustrated in Figure 25C. The freedom to move is primarily in the direction of molecular alignment, as sideways motion quickly results in collision with a neighbour. Layering the grains, as illustrated in Figure 25B, enhances sideways motion. This produces the smectic phase. In the smectic phase some molecules have ample free volume to move in, while others are tightly packed. The lowest-energy arrangement shares the free volume equitably among molecules. Each molecular environment matches all others, and the structure is a crystal like that illustrated in Figure 25A.

There is a great variety of liquid crystalline structures known in addition to those described so far. Table 2 relates some of the chief structures according to their degree and type of order. The smectic-C phase and those listed below it have molecules tilted with respect to the layers. Continuous in-plane rotational symmetry, present within smectic-A layers, is broken in the hexatic-B phase, but a proliferation of dislocations maintains continuous translational symmetry within its layers. A similar relationship holds between smectic-C and smectic-F. Crystal-B and crystal-G have molecular positions on regular crystal lattice sites, with long axes of molecules (directors) aligned, but allow rotation of molecules about their directors. These are the so-called plastic crystals. Many interesting liquid crystal phases are not listed in this table, including the discotic phase, consisting of disk-shaped molecules, and the columnar phases, in which translational symmetry is broken in not one but two spatial directions, leaving liquidlike order only along columns. The degree of order increases from the top to the bottom of the table. In general, phases from the top of the table are expected at high temperatures, and phases from the bottom at low temperatures.

Rice-grain analogy

Symmetry of liquids

Liquid crystal structures

Table 2: Selected Phases Characteristic of Liquid-Crystal-Forming Molecules

phase		order
Isotropic liquid Nematic		full continuous translational and rotational symmetry molecular orientation breaks rotational symmetry
untilted	tilted	
Smectic-A Hexatic-B Crystal-B Crystal-E	Smectic-C Smectic-F Crystal-G Crystal-H	layering breaks translational symmetry; smectic-C molecules are tilted bond orientational order breaks rotational symmetry within layers crystallization breaks translational symmetry within layers; molecules may rotate about their long axis molecular rotation freezes out

Structure of soap molecules

Liquid-crystal-forming compounds are widespread and quite diverse. Soap (see Figure 26) can form a type of smectic known as a lamellar phase, also called neat soap. In this case it is important to recognize that soap molecules have a dual chemical nature. One end of the molecule (the hydrocarbon tail) is attracted to oil, while the other end (the polar head) attaches itself to water. When soap is placed in water, the hydrocarbon tails cluster together, while the polar heads adjoin the water. Small numbers of soap molecules form spherical or rodlike micelles (Figure 26B), which float freely in the water, while concentrated solutions create bilayers (Figure 26C), which stack along some direction just like smectic layers. Indeed, the name smectic is derived from the Greek word for soap. The slippery feeling caused by soap reflects the ease with which the layers slide across one another.

Biologically important liquid crystals

Many biological materials form liquid crystals. Myelin, a fatty material extracted from nerve cells, was the first intensively studied liquid crystal. The tobacco mosaic virus, with its rodlike shape, forms a nematic phase. In cholesterol the nematic phase is modified to a cholesteric phase characterized by continuous rotation of the direction of molecular alignment. An intrinsic twist of the cholesterol molecule, rather like the twist of the threads of a screw, causes this rotation. Since the molecular orientation rotates steadily, there is a characteristic distance after which the orientation repeats itself. This distance is frequently comparable to the wavelength of visible light, so brilliant colour effects result from the diffraction of light by these materials.

Perhaps the first description of a liquid crystal occurred in the story *The Narrative of Arthur Gordon Pym*, by Edgar Allan Poe:

I am at a loss to give a distinct idea of the nature of this liquid, and cannot do so without many words. Although it flowed with rapidity in all declivities where common water would do so, yet never, except when falling in a cascade, had it the customary appearance of limpidity. . . . At first sight, and especially in cases where little declivity was found, it bore resemblance, as regards consistency, to a thick infusion of gum Arabic in common water. But this was only the least remarkable of its extraordinary qualities. It was *not* colourless, nor was it of any one uniform colour—presenting to the eye, as it flowed, every possible shade of purple, like the hues of a changeable silk. . . . Upon collecting a basinful, and allowing

it to settle thoroughly, we perceived that the whole mass of liquid was made up of a number of distinct veins, each of a distinct hue; that these veins did not commingle; and that their cohesion was perfect in regard to their own particles among themselves, and imperfect in regard to neighbouring veins. Upon passing the blade of a knife athwart the veins, the water closed over it immediately, as with us, and also, in withdrawing it, all traces of the passage of the knife were instantly obliterated. If, however, the blade was passed down accurately between two veins, a perfect separation was effected, which the power of cohesion did not immediately rectify.

The liquid described in this passage is human blood. In its usual state within the human body, blood is an ordinary disordered isotropic fluid. The disklike shape of red blood cells, however, favours liquid crystallinity at certain concentrations and temperatures.

OPTICAL PROPERTIES

An understanding of the principal technological applications of liquid crystals requires a knowledge of their optical properties. Liquid crystals alter the polarization of light passing through them. Light waves are actually waves in electric and magnetic fields. The direction of the electric field is the polarization of the light wave. A polarizing filter selects a single component of polarized light to pass through while absorbing all other components of incoming waves. If a second polarizing filter is placed above the first but with its polarization axis rotated by 90° , no light can pass through because the polarization passed by the first filter is precisely the polarization blocked by the second filter. When optically active materials, such as liquid crystals, are placed between polarizing filters crossed in this manner, some light may get through, because the intervening material changes the polarization of the light. If the nematic director is not aligned with either of the polarizing filters, polarized light passing through the first filter becomes partially polarized along the nematic director. This component of light in turn possesses a component aligned with the top polarizing filter, so a fraction of the incoming light passes through the entire assembly. The amount of light passing through is largest when the nematic director is positioned at a 45° angle from both filters. The light is fully blocked when the director lies parallel to one filter or the other.

During the last decades of the 19th century, pioneering

From G.H. Brown, J.D. Doane, and V.D. Neff, *A Review of the Structure and Physical Properties of Liquid Crystals*, CRC Press, Inc. 1971, reproduced by permission of the *Journal of the Society of Cosmetic Chemists*

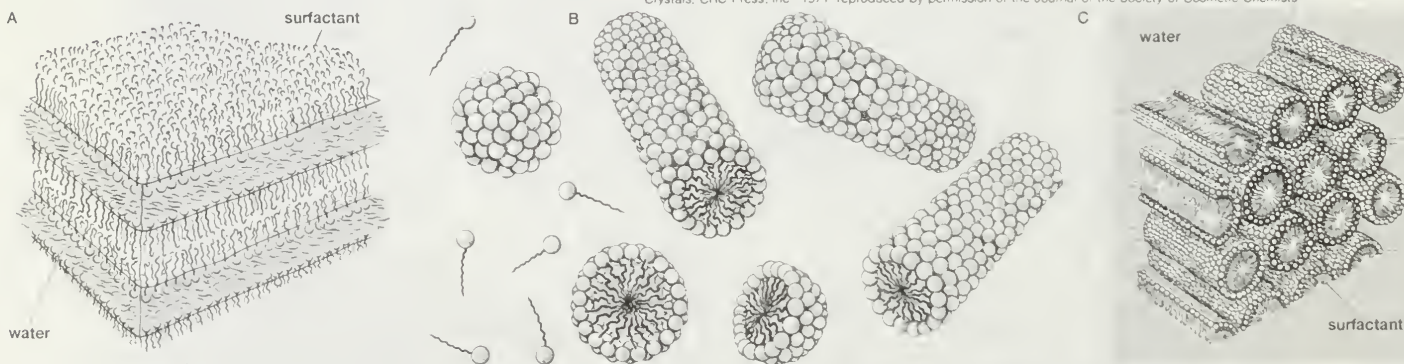


Figure 26: Structures of soap in water.

(A) The smectic phase, also called the neat phase. (B) Spherical and rodlike micelles formed by soap molecules; they float freely in dilute solutions. (C) Bilayer packing of rodlike micelles in a concentrated solution.

investigators of liquid crystals, such as the German physicist Otto Lehmann and the Austrian botanist Friedrich Reinitzer, equipped ordinary microscopes with pairs of polarizing filters. Typical microscope images of nematic and smectic phases taken through crossed polarizers are shown in Figure 27. Spatial variation in the alignment of the nematic director causes spatial variation in light intensity. Since the nematic is defined by having all di-

rectors nearly parallel to one another, the images arise from defects in the nematic structure. Figure 27 (bottom) illustrates a manner in which the directors may rotate or bend around defect lines. The resulting threadlike images inspired the name nematic, which is based on the Greek word for thread. The layered smectic structure causes layering of defects in Figure 27 (centre).

Nonuniformity in director alignment may be induced

(Top and bottom) From Jurgen Nehring and Alfred Saupe, *Journal of The Chemical Society, Faraday Transactions II*, 1972, vol. 68, 1-15. © copyright 1972 by The Chemical Society, London, (centre) Dietrich Demus

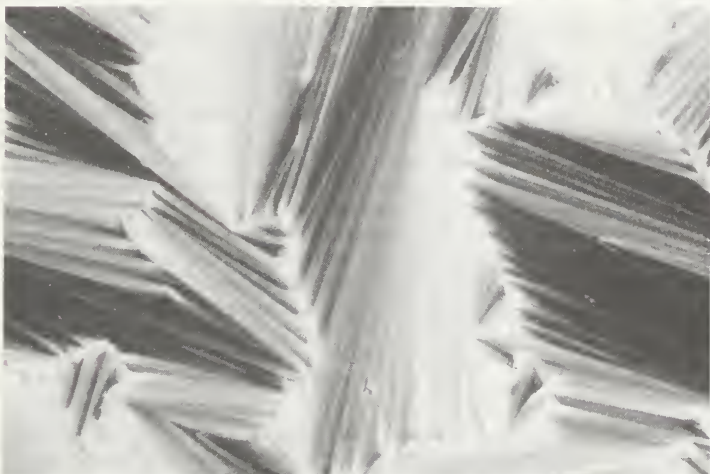
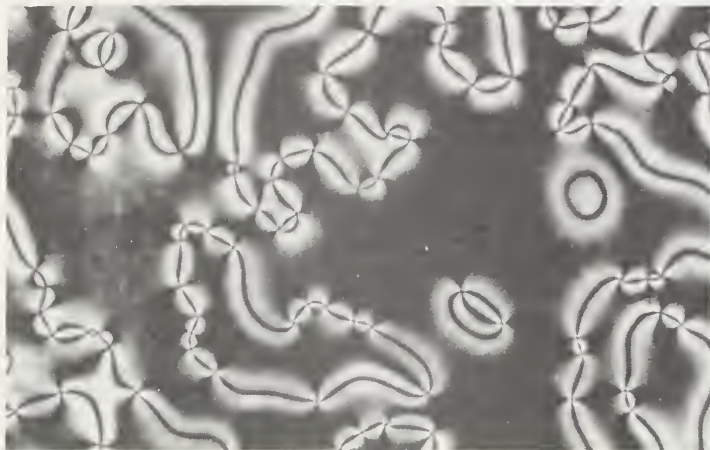


Figure 27: (Top) Nematic and (centre) smectic liquid crystals viewed through microscopes equipped with crossed polarizers. (Bottom) Spatial variation of the director, causing the threadlike image seen in the photograph at top.

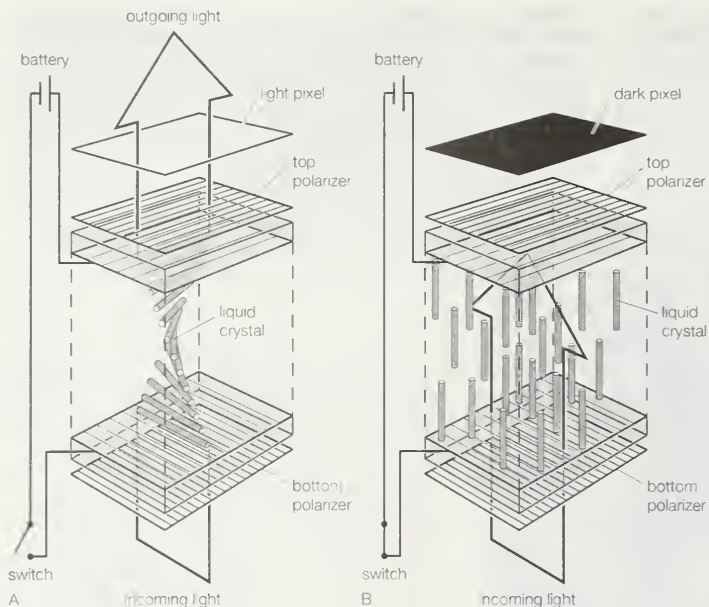


Figure 28: A twisted-nematic cell. (A) The assembly is transparent to light in the absence of an electric field. (B) An applied field destroys the twist of the nematic, rendering the assembly opaque.

From J. Funtschilling, "Liquid Crystals and Liquid Crystal Displays," *Condensed Matter News*, vol. 1, no. 1, 1991. Gordon and Breach Science Publishers.

artificially. The surfaces of a glass container can be coated with a material that, when rubbed in the proper direction, forces the director to lie perpendicular or parallel to the wall adjacent to a nematic liquid crystal. The orientation forced by one wall need not be consistent with that forced by another wall; this situation causes the director orientation to vary in between the walls. The nematic must compromise its preference for all directors to be parallel to one another with the inconsistent orienting forces of the container walls. In doing so, the liquid crystal may take on a twisted alignment across the container (see Figure 28A.) Electric or magnetic fields provide an alternate means of influencing the orientation of the nematic directors. Molecules may prefer to align so that their director is, say, parallel to an applied electric field.

The twisted-nematic cell

Optical behaviour and orienting fields underlie the important contemporary use of liquid crystals as optoelectronic displays. Consider, for example, the twisted-nematic cell shown in Figure 28A. The polarizer surfaces are coated and rubbed so that the nematic will align with the polarizing axis. The two polarizers are crossed, forcing the nematic to rotate between them. The rotation is slow and smooth, assuming a 90° twist across the cell. Light passing through the first polarizer is aligned with the bottom of the nematic layer. As the nematic twists, it rotates the polarization of the light so that, as the light leaves the top of the nematic layer, its polarization is rotated by 90° from that at the bottom. The new polarization is just

right for passing through the top filter, and so light travels unhindered through the assembly.

If an electric field is applied in the direction of light propagation, the liquid crystal directors align with the orienting field, so they are no longer parallel to the light passing through the bottom polarizer (Figure 28B). They are no longer capable of rotating this polarization through the 90° needed to allow the light to emerge from the top polarizer. Although this assembly is transparent when no field is applied, it becomes opaque when the field is present. A grid of such assemblies placed side by side may be used to display images. If one turns on the electric field attached to the parts of the grid that lie where the image is to appear, these points will turn black while the remaining points of the grid stay white. The resulting patchwork of dark and light creates the image on the display. In a wristwatch, calculator, or computer these may be simply numbers or letters, and in a television the images may be detailed pictures. Switching the electric fields on or off will cause the picture to move, just as ordinary television pictures display an ever-changing stream of electrically encoded images. (M.W.)

Amorphous solids

Solids and liquids are both forms of condensed matter; both are composed of atoms in close proximity to each other. But their properties are, of course, enormously

From R. Zallen, *The Physics of Amorphous Solids*, copyright © 1983 John Wiley & Sons, Inc. reprinted by permission of John Wiley & Sons, Inc.

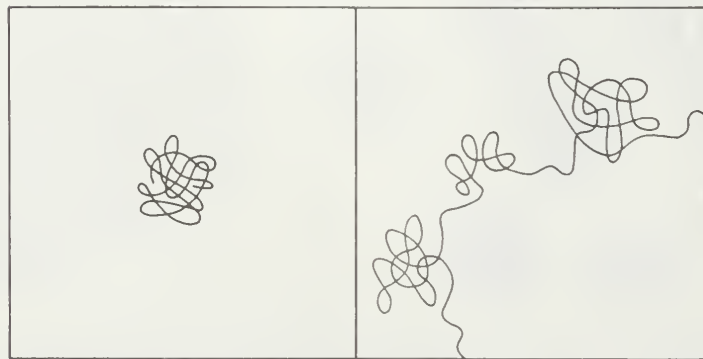


Figure 29: The state of atomic motion (left) in a solid and (right) in a liquid.

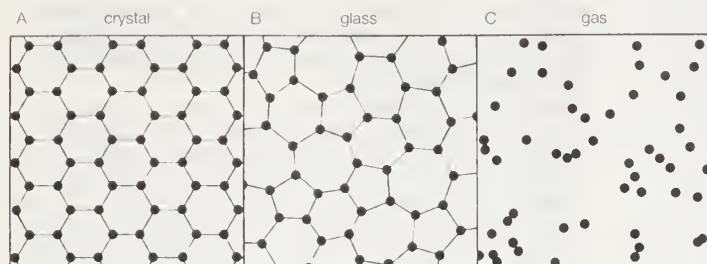


Figure 30: The atomic arrangements in (A) a crystalline solid, (B) an amorphous solid, and (C) a gas.

From R. Zallen, *The Physics of Amorphous Solids*, copyright © 1983 John Wiley & Sons, Inc., reprinted by permission of John Wiley & Sons, Inc.

Defining characteristics of solids

different. While a solid material has both a well-defined volume and a well-defined shape, a liquid has a well-defined volume but a shape that depends on the shape of the container. Stated differently, a solid exhibits resistance to shear stress while a liquid does not. Externally applied forces can twist or bend or distort a solid's shape, but (provided the forces have not exceeded the solid's elastic limit) it "springs back" to its original shape when the forces are removed. A liquid flows under the action of an external force; it does not hold its shape. These macroscopic characteristics constitute the essential distinctions: a liquid flows, lacks a definite shape (though its volume is definite), and cannot withstand a shear stress; a solid does not flow, has a definite shape, and exhibits elastic stiffness against shear stress.

On an atomic level, these macroscopic distinctions arise from a basic difference in the nature of the atomic motion. Figure 29 contains schematic representations of atomic movements in a liquid and a solid. Atoms in a solid are not mobile. Each atom stays close to one point in space, although the atom is not stationary but instead oscillates rapidly about this fixed point (the higher the temperature, the faster it oscillates). The fixed point can be viewed as a time-averaged centre of gravity of the rapidly jiggling atom. The spatial arrangement of these fixed points constitutes the solid's durable atomic-scale structure. In contrast, a liquid possesses no enduring arrangement of atoms. Atoms in a liquid are mobile and continually wander throughout the material.

DISTINCTION BETWEEN CRYSTALLINE AND AMORPHOUS SOLIDS

There are two main classes of solids: crystalline and amorphous. What distinguishes them from one another is the nature of their atomic-scale structure. The essential differences are displayed in Figure 30. The salient features of the atomic arrangements in amorphous solids (also called glasses), as opposed to crystals, are illustrated in the figure for two-dimensional structures; the key points carry over to the actual three-dimensional structures of real materials. Also included in the figure, as a reference material, is a sketch of the atomic arrangement in a gas. For the sketches representing crystal (A) and glass (B) structures, the solid dots denote the fixed points about which the atoms oscillate; for the gas (C), the dots denote a snapshot of one configuration of instantaneous atomic positions.

Atomic positions in a crystal exhibit a property called long-range order or translational periodicity; positions repeat in space in a regular array, as in Figure 30A. In an amorphous solid, translational periodicity is absent. As indicated in Figure 30B, there is no long-range order. The atoms are not randomly distributed in space, however, as they are in the gas in Figure 30C. In the glass example illustrated in the figure, each atom has three nearest-neighbour atoms at the same distance (called the chemical bond length) from it, just as in the corresponding crystal. All solids, both crystalline and amorphous, exhibit short-range (atomic-scale) order. (Thus, the term amorphous, literally "without form or structure," is actually a misnomer in the context of the standard expression amorphous solid.) The well-defined short-range order is a consequence of the chemical bonding between atoms, which is responsible for holding the solid together.

In addition to the terms amorphous solid and glass,

other terms in use include noncrystalline solid and vitreous solid. Amorphous solid and noncrystalline solid are more general terms, while glass and vitreous solid have historically been reserved for an amorphous solid prepared by rapid cooling (quenching) of a melt—as in scenario 2 of Figure 31.

Figure 31, which should be read from right to left, indicates the two types of scenarios that can occur when cooling causes a given number of atoms to condense from the gas phase into the liquid phase and then into the solid phase. Temperature is plotted horizontally, while the volume occupied by the material is plotted vertically. The temperature T_b is the boiling point, T_f is the freezing (or melting) point, and T_g is the glass transition temperature. In scenario 1 the liquid freezes at T_f into a crystalline solid, with an abrupt discontinuity in volume. When cooling occurs slowly, this is usually what happens. At sufficiently high cooling rates, however, most materials display a different behaviour and follow route 2 to the solid state. T_f is bypassed, and the liquid state persists until the lower temperature T_g is reached and the second solidification scenario is realized. In a narrow temperature range near T_g , the glass transition occurs: the liquid freezes into an amorphous solid with no abrupt discontinuity in volume.

The glass transition temperature T_g is not as sharply defined as T_f ; T_g shifts downward slightly when the cooling rate is reduced. The reason for this phenomenon is the steep temperature dependence of the molecular response time, which is crudely indicated by the order-of-magnitude values shown along the top scale of Figure 31. When the temperature is lowered below T_g , the response time for molecular rearrangement becomes much larger than experimentally accessible times, so that liquidlike mobility (Figure 29, right) disappears and the atomic configuration becomes frozen into a set of fixed positions to which the atoms are tied (Figures 29, left, and 30B).

Some textbooks erroneously describe glasses as undercooled viscous liquids, but this is actually incorrect. Along the section of route 2 labeled liquid in Figure 31, it is the

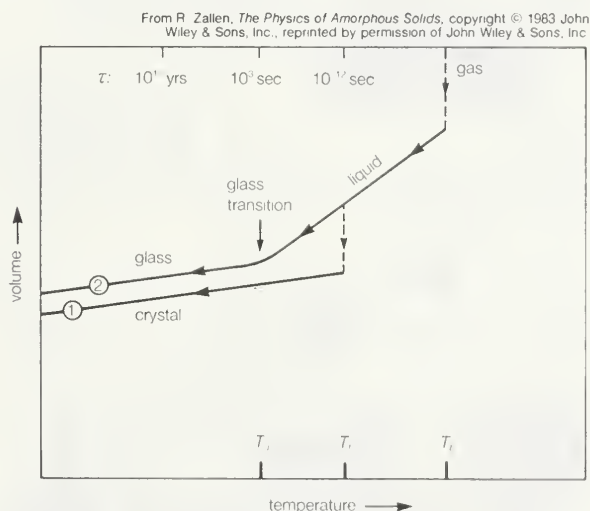


Figure 31: The two general cooling paths by which a group of atoms can condense. Route 1 is the path to the crystalline state; route 2 is the rapid-quench path to the amorphous solid state.

Short-range order

Table 3: Bonding Types and Glass Transition Temperatures of Representative Amorphous Solids

glass	bonding	T_g (K)
SiO ₂	covalent	1,430
GeO ₂	covalent	820
Si, Ge	covalent	—
Pd _{0.4} Ni _{0.4} P _{0.2}	metallic	580
BeF ₂	ionic	570
As ₂ S ₃	covalent	470
Polystyrene	polymeric	370
Se	polymeric	310
Au _{0.8} Si _{0.2}	metallic	290
H ₂ O	hydrogen-bonded	140
C ₂ H ₅ OH	hydrogen-bonded	90
Isopentane	van der Waals	65
Fe, Co, Bi	metallic	—

portion lying between T_f and T_g that is correctly associated with the description of the material as an undercooled liquid (undercooled meaning that its temperature is below T_f). But below T_g in the glass phase, it is a bona fide solid (exhibiting such properties as elastic stiffness against shear). The low slopes of the crystal and glass line segments of Figure 31 in comparison with the high slope of the liquid section reflect the fact that the coefficient of thermal expansion of a solid is small in comparison with that of the liquid.

PREPARATION OF AMORPHOUS SOLIDS

It was once thought that relatively few materials could be prepared as amorphous solids, and such materials (notably, oxide glasses and organic polymers) were called glass-forming solids. It is now known that the amorphous solid state is almost a universal property of condensable matter. Table 3 presents a list of amorphous solids in which every class of chemical bonding type is represented. The glass transition temperatures span a wide range.

Glass formation is a matter of bypassing crystallization. The channel to the crystalline state is evaded by quickly crossing the temperature interval between T_f and T_g . Nearly all materials can, if cooled quickly enough, be prepared as amorphous solids. The definition of "quickly enough" varies enormously from material to material. Four techniques for preparing amorphous solids are illustrated in Figure 32. These techniques are not fundamentally different from those used for preparing crystalline

solids: the key is simply to quench the sample quickly enough to form the glass, rather than slowly enough to form the crystal. The quench rate increases greatly from left to right in the figure.

Preparation of metallic glasses requires a quite rapid quench. The technique shown in Figure 32C, called splat quenching, can quench a droplet of a molten metal roughly 1,000° C in one millisecond, producing a thin film of metal that is an amorphous solid. In enormous contrast to this, the silicate glass that forms the rigid ribbed disk of the Hale telescope of the Palomar Observatory near San Diego, Calif., was prepared by cooling (over a comparable temperature drop) during a time interval of eight months. The great difference in the quench rates needed for arriving at the amorphous solid state (the quench rates here differ by a factor of 3×10^{10}) is a dramatic demonstration of the difference in the glass-forming tendency of silicate glasses (very high) and metallic glasses (very low).

The required quench rate for glass formation can vary significantly within a family of related materials that differ from one another in chemical composition. Figure 33 illustrates a representative behaviour for a binary (two-component) system, gold-silicon. Here x specifies the fraction of atoms that are silicon atoms, and Au_{1-x}Si_x denotes a particular material in this family of materials. (Au is the chemical symbol for gold, Si is the symbol for silicon, and, for example, Au_{0.8}Si_{0.2} denotes a material containing 20 percent silicon atoms and 80 percent gold atoms.) The solid curve labeled T_f shows the composition dependence of the freezing point; above this line the liquid phase is the stable form. There is a deep cusp near the composition $x = 0.2$. Near this special composition, as at a in the figure, a liquid is much more readily quenched than is a liquid at a distant composition such as b . To reach the glass phase, the liquid must be cooled from above T_f to below T_g without crystallizing. Throughout the temperature interval from T_f down to the glass transition temperature T_g , the liquid is at risk vis-à-vis crystallization. Since this dangerous interval is much longer at b than at a , a faster quench rate is needed for glass formation at b than at a .

Diagrams similar to (though slightly more complicated than) Figure 33 exist for many binary systems. For example, in the oxide system CaO-Al₂O₃, in which the two end-member compositions ($x = 0$ and $x = 1$) correspond to pure calcium oxide (CaO) and pure aluminum oxide

Splat
quenching

From R. Zallen, *The Physics of Amorphous Solids*, copyright © 1983 John Wiley & Sons, Inc., reprinted by permission of John Wiley & Sons, Inc.

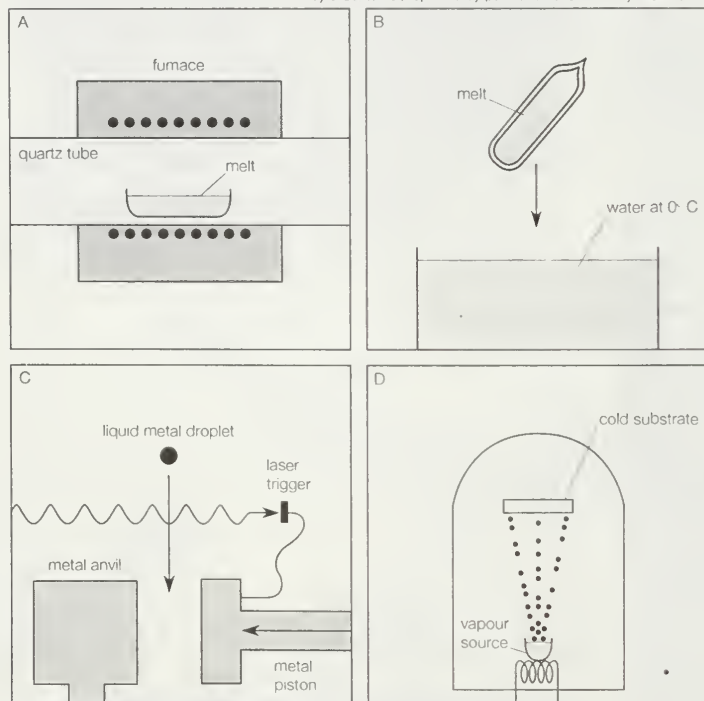


Figure 32: Four methods for preparing amorphous solids. (A) Slow cooling, (B) moderate quenching, (C) rapid splat quenching, and (D) condensation from the gas phase.

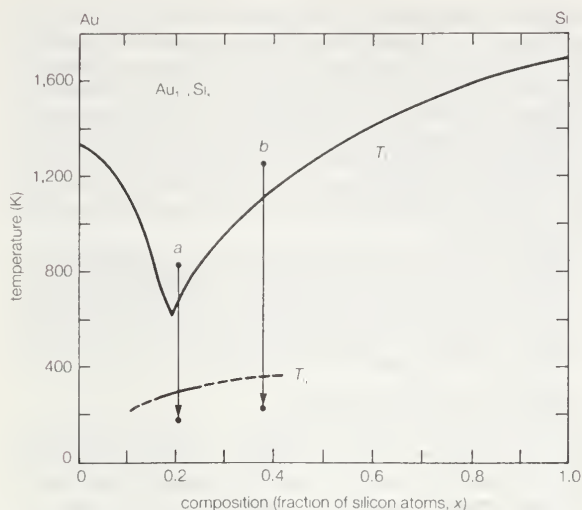


Figure 33: Glass formation in the gold-silicon system. Two quenches from the liquid state are shown. Glasses can be prepared more easily with quench a than with quench b, because the latter requires the liquid to cross a larger temperature interval between the freezing temperature T_f and the glass transition temperature T_g ; this is the temperature region in which crystallization is prone to occur.

Based on R. Zallen, *The Physics of Amorphous Solids*, copyright © 1983 John Wiley & Sons, Inc., reprinted by permission of John Wiley & Sons, Inc. T_f curve from B. Predel and H. Bankstahl, *Journal of the Less-Common Metals*, no. 43, 1975, Elsevier Sequoia, publisher. T_g curve from H.S. Chen and D. Turnbull, *The Journal of Chemical Physics*, no. 48, 1968, published by the American Institute of Physics.

(Al_2O_3), there is a deep minimum in the T_f -versus- x curve near the middle of the composition range. Although neither calcium oxide nor aluminum oxide readily forms a glass, glasses are easily formed from mixed compositions; for reasons related to this, many oxide glasses have complex chemical compositions.

In the gold-silicon system of Figure 33, at compositions far from the cusp, glasses cannot be formed by melt quenching—even by the rapid splat-quench technique of Figure 32. (This is the reason that the T_g curve of Figure 33 spans only compositions near the cusp.) Amorphous solids can still be prepared by dispensing with the liquid phase completely and constructing a thin solid film in atom-by-atom fashion from the gas phase. Figure 32D shows the simplest of these vapour-condensation techniques. A vapour stream, formed within a vacuum chamber by thermal evaporation of a sample of the material to be deposited, impinges on the surface of a cold substrate. The atoms condense on the cold surface and, under a range of conditions (usually a high rate of deposition and a low substrate temperature), an amorphous solid is formed as a thin film. Pure silicon can be prepared as an amorphous solid in this manner. Variations of the method include using an electron beam to vapourize the source or using the plasma-induced decomposition of a molecular species. The latter technique is used to deposit amorphous silicon from gaseous silane (SiH_4). Among the amorphous solids listed in Table 3, those that normally require vapour-condensation methods for their preparation are silicon (Si), germanium (Ge), water (H_2O), and the elemental metallic glasses iron (Fe), cobalt (Co), and bismuth (Bi).

Numerous other methods exist for preparing amorphous solids, and new methods are continually invented. In melt spinning, a jet of molten metal is propelled against the moving surface of a cold, rotating copper cylinder. A solid film of metallic glass is spun off as a continuous ribbon at a speed that can exceed a kilometre per minute. In laser glazing, a brief intense laser pulse melts a tiny spot, which is swiftly quenched by the surrounding material into a glass. In sol-gel synthesis, small molecules in a liquid solution chemically link up with each other, forming a disordered network. It is possible to take a crystalline solid and convert it into an amorphous solid by bombarding it with high-kinetic-energy ions. Under certain conditions of composition and temperature, interdiffusion (mixing on an atomic scale) between crystalline layers can produce an amorphous phase. Pyrolysis and electrolysis are other methods that can be used.

ATOMIC-SCALE STRUCTURE

The absence of long-range order is the defining characteristic of the atomic arrangement in amorphous solids. However, because of the absence in glasses of long parallel rows and flat parallel planes of atoms, it is extremely difficult to determine details of the atomic arrangement with the structure-probing techniques (such as X-ray diffraction) that are so successful for crystals. For glasses the information obtained from such structure-probing experiments is contained in a curve called the radial distribution function (RDF).

The radial distribution function

Figure 34 shows a comparison of the experimentally determined RDFs of the crystalline and amorphous forms of germanium, an elemental semiconductor similar to silicon. The heavy curve labeled a-Ge corresponds to amorphous germanium; the light curve labeled c-Ge corresponds to crystalline germanium. The significance of the RDF is that it gives the probability of neighbouring atoms being located at various distances from an average atom. The horizontal axis in the figure specifies the distance from a given atom; the vertical axis is proportional to the average number of atoms found at each distance. (The distance scale is expressed in angstrom units; one angstrom equals 10^{-8} centimetre.) The curve for crystalline germanium displays sharp peaks over the full range shown, corresponding to well-defined shells of neighbouring atoms at specific distances, which arise from the long-range regularity of the crystal's atomic arrangement. Amorphous germanium exhibits a close-in sharp peak corresponding to the nearest-neighbour atoms (there are four nearest neighbours in both c-Ge and a-Ge), but at larger distances the undulations in the RDF curve become washed out owing to the absence of long-range order. The first, sharp, nearest-neighbour peak in a-Ge is identical to the corresponding peak in c-Ge, showing that the short-range order in the amorphous form of solid germanium is as well-defined as it is in the crystalline form.

The detailed shape of the a-Ge RDF curve of Figure 34 is the input used in the difficult task of developing a model for the atomic arrangement in amorphous germanium. The normal procedure is to construct a model of the structure and then to calculate from the model's atomic positions a theoretical RDF curve. This calculated RDF is then compared to the experimental curve (which provides the definitive test of the validity of the model). Computer-assisted refinements are then made in the model in order to improve the agreement between the model-dependent theoretical RDF and the experimentally observed RDF. This program has been successfully carried out for many amorphous solids, so there is now much that is known about their atomic-scale structure. In contrast to the complete information available for crystals, however, the structural knowledge of glasses still contains gaps.

Amorphous solids, like crystalline solids, exhibit a wide

After R.J. Temkin, W. Paul, and G.A.N. Connell, *Advances in Physics*, no. 22, 1973, Taylor and Francis Ltd., publisher, in R. Zallen, *The Physics of Amorphous Solids*, copyright © 1983 John Wiley & Sons, Inc., reprinted by permission of John Wiley & Sons, Inc.

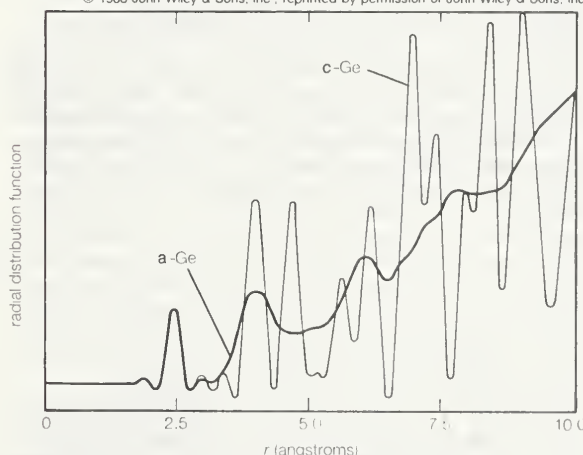


Figure 34: Comparison of the atomic radial distribution functions of crystalline (c-Ge) and amorphous (a-Ge) germanium. The value of the function at each distance r from a given atom is proportional to the number of atoms found at that distance.

Models of
atomic-
scale
structure

variety of atomic-scale structures. Most of these can be recognized as falling within one or another of three broad classes of structure associated with the following models: (1) the continuous random-network model, applicable to covalently bonded glasses, such as amorphous silicon and the oxide glasses, (2) the random-coil model, applicable to the many polymer-chain organic glasses, such as polystyrene, and (3) the random close-packing model, applicable to metallic glasses, such as $Au_{0.8}Si_{0.2}$ gold-silicon. These are the names in conventional use for the models. Although each of them contains the word random, the well-defined short-range order means that they are not random in the sense that the gas structure of Figure 30C is random.

An illustration of the continuous random-network model is shown in Figure 35A and of the random-coil model in Figure 35B. Figure 35A reproduces a famous diagram published by W.H. Zachariasen in 1932. It is for a hypothetical two-dimensional A_2B_3 glass in which every A atom is bonded to three B atoms and every B atom to two A atoms. This picture bears a reasonable resemblance to current models for the arsenic chalcogenide glasses As_2S_3 and As_2Se_3 . (Sulfur, S, and selenium, Se, belong to the group of elements called chalcogens.) The model was introduced as a schematic analogue for the network structure of the oxide glasses. The prototypical oxide glass is amorphous SiO_2 , or silica glass. (Quartz, which is present in sand, is a crystalline form of SiO_2 .) In amorphous SiO_2 , each silicon atom is bonded to four oxygen atoms, and each oxygen atom is bonded to two silicon atoms. This structure is difficult to represent in a two-dimensional picture, but Figure 35A is a useful analogue with the hollow circles representing oxygen atoms and the small solid dots representing silicon atoms. The fourth bond originating at each silicon can be imagined to be out of the plane of the diagram.

From R. Zallen, *The Physics of Amorphous Solids*, copyright © 1983 John Wiley & Sons, Inc., reprinted by permission of John Wiley & Sons, Inc., (A) from W.H. Zachariasen, *Journal of the American Chemical Society*, no. 54 1932

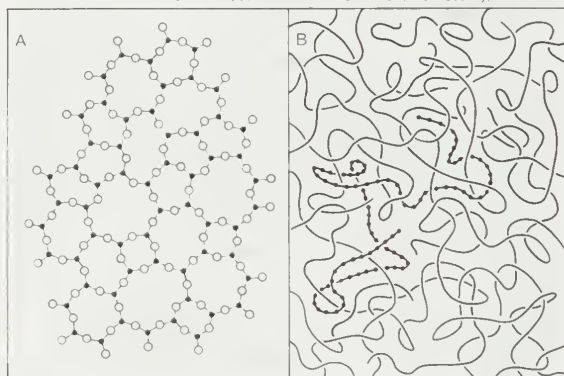


Figure 35: Two basic models for the atomic-scale structure of amorphous solids.

(A) The continuous random-network model for network glasses, and (B) the random-coil model for polymeric glasses.

The network structure shown in Figure 35A clearly demonstrates how short-range order (note the triangle of neighbours surrounding each solid dot) is compatible with the absence of long-range order. At the bridging oxygen atoms, the bond angles have some flexibility, so it is easy to continue the network. Common oxide glasses are chemically more complex than SiO_2 , as discussed in the next section. Chemical species such as phosphorus and germanium, which (like silicon) enter into the structure of the network by forming strong chemical bonds with oxygen atoms, are called network formers. Chemical species such as sodium and calcium, which do not bond directly to the network but which simply sit (in ionic form) within its interstitial holes, are called network modifiers.

A large fraction of the everyday materials called plastics are amorphous solids composed of long-chain molecules known as polymers. Each polymer chain has a backbone consisting of a string of many (up to roughly 100,000) carbon atoms bonded to each other. These organic polymeric glasses are present in innumerable familiar molded products (e.g., pens, tires, toys, appliance bodies, building ma-

terials, and automobile and airplane parts). The random-coil model of Figure 35B, first proposed in 1949 by P.J. Flory (who later received a Nobel Prize in Chemistry for his pioneering work on polymers), is the established structural model for this important class of amorphous solids. As schematically sketched in the figure, the structure consists of intermeshed, entangled polymer chains. The chain configurations are well-defined, statistically, by a mathematical trajectory called a three-dimensional random walk.

The third important structural model, the random close-packing model for metallic glasses, is difficult to illustrate with a simple diagram. Roughly speaking, it is similar to the structure that arises when a bunch of marbles are swiftly scrunched together in a paper bag. (R.Z.)

PROPERTIES OF OXIDE GLASSES

The wide range of the properties of glasses depends on their composition, and special effects result from the presence of various modifying agents in certain basic glass-forming materials (see above *Atomic-scale structure*).

One of the most important glass formers is silica (SiO_2). Pure crystalline silica melts at $1,710^\circ C$. In pure form, silica glass exhibits such properties as low thermal expansion, high softening temperatures, and excellent chemical and electrical resistance. In pure form it is relatively transparent over a wide range of wavelengths to visible and ultraviolet light and to ultrasonic waves.

The high viscosity (see below) and melting temperature of silica glass are affected by the presence or absence of other materials. For example, if certain materials called fluxes are added, the most important being soda (Na_2O), both viscosity and melting temperature can be reduced. If too much soda is added, the resulting glass is readily attacked by water, but, if there are suitable amounts of stabilizing oxides, such as lime (CaO) and magnesia (MgO), the glass becomes more durable. Most commercial glass has a soda-lime-silica composition and is produced in vast quantities for plate and sheet glass, containers, and lightbulbs.

In soda-lime-silica glasses, if lime is replaced by lead oxide (PbO) and if potash (K_2O) is used as a partial replacement for soda, lead-alkali-silicate glasses result that have lower softening points than lime glasses. The refractive indices, dispersive powers, and electrical resistance of these glasses are generally much greater than those of soda-lime-silica glasses.

Boric oxide (B_2O_3), itself a glass former, acts as a flux (i.e., lowers the working temperature) when present in silica and forms borosilicate glass, and the substitution of small percentages of alkali and alumina increases the chemical stability. It also exhibits low thermal expansion, high dielectric strength, and high softening temperature.

Aluminosilicate glasses find applications similar to those of borosilicates, but the former can stand higher operating temperatures; glasses with relatively high alumina contents and no boric oxide are exceptionally resistant to alkalis.

The above glasses all have silica as the glass former. With other glass formers, glasses have special properties. For example, if boric oxide is present, X rays are transmitted and rare-earth glasses will exhibit low dispersion and a high refractive index. Phosphate glasses (used as optical glasses) based on phosphorus pentoxide (P_2O_5) are highly resistant to hydrofluoric acid and act as efficient heat absorbers when iron oxide is added. Table 4 gives the compositions and physical properties of some typical commercial oxide glasses of the types described. (R.W.D./R.Z.)

PROPERTIES AND APPLICATIONS OF AMORPHOUS SOLIDS

The following sections discuss technological applications of amorphous solids in connection with the properties that make those applications possible. It is important to understand that, although differences do exist between the properties of amorphous and crystalline solids, it is nevertheless broadly true that amorphous solids exhibit essentially the full range of properties and phenomena exhibited by crystalline solids. There are amorphous-solid metals, semiconductors, and insulators; there are transparent glasses and opaque glasses; and there are superconducting amorphous solids and ferromagnetic amorphous solids.

Silica glass
formers

Nonsilica
glass
formers

Table 4: Characteristics of Oxide Glasses

	fused silica	soda-lime-silica	borosilicate	aluminosilicate	lead
Approximate composition*	SiO ₂ 99.9%	SiO ₂ 73%	SiO ₂ 81%	SiO ₂ 62%	SiO ₂ 56%
	H ₂ O 0.1%	Al ₂ O ₃ 1%	Al ₂ O ₃ 2%	Al ₂ O ₃ 17%	Al ₂ O ₃ 2%
		Na ₂ O 17%	B ₂ O ₃ 13%	B ₂ O ₃ 5%	Na ₂ O 4%
		MgO 4%	Na ₂ O 4%	Na ₂ O 1%	K ₂ O 9%
		CaO 5%		MgO 7%	PbO 29%
			CaO 8%		
Coefficient of thermal expansion (linear expansion per °C × 10 ⁷)	5.5	93	33	42	89
Strain point † (°C; viscosity about 10 ^{14.5} poise)	990	470	515	670	395
Annealing point ‡ (°C; viscosity about 10 ¹³ poise)	1,050	510	565	715	435
Softening point § (°C; viscosity about 10 ^{7.65} poise)	1,580	695	820	915	630
Young's modulus (lbs per sq in. × 10 ⁶)	10.5	10	9.1	12.7	8.6
Refractive index (for sodium D line)	1.459	1.512	1.474	1.530	1.560
Dielectric constant (at 10 ⁶ cycles per second and 20° C)	3.8	7.2	4.6	7.2	6.7
Density (g/cm ³)	2.20	2.47	2.23	2.52	3.05

*These compositions are typical of the various glass types. †Temperature at which internal stresses are reduced significantly over a few hours. ‡Temperature at which internal stresses are reduced significantly over a few minutes. §Temperature at which glass will rapidly deform under its own weight. ||The strain point and annealing point roughly define the annealing range.

Some of the general differences between the properties of crystals and glasses, in addition to the fundamental one of the glass transition (as discussed above in connection with Figure 31 and also below with regard to its value in technological settings), are noted here. The atomic-scale disorder present in a metallic glass causes its electrical conductivity to be lower than the conductivity of the corresponding crystalline metal, because the structural disorder impedes the motion of the mobile electrons that make up the electrical current. (This lower electrical conductivity for the amorphous metal can be an advantage in some situations, as discussed below in the section *Magnetic glasses*.) For a similar reason, the thermal conductivity of an insulating glass is lower than that of the corresponding crystalline insulator; glasses thus make good thermal insulators. Crystals and glasses also differ systematically in their optical spectra, which are the curves that describe the wavelength dependence of the degree to which the solid absorbs infrared, visible, or ultraviolet light. Although the overall spectra are often similar, crystal spectra typically exhibit sharp peaks and other features that specifically arise as a consequence of the long-range order of the crystal's atomic-scale structure. These sharp features are absent in the optical spectra of amorphous solids.

The continuous liquid-to-solid transition near T_g , the glass transition, has a profound significance in connection with classical applications of glasses. While crystallization abruptly transforms a mobile, low-viscosity liquid to a crystalline solid at T_f , near T_g the liquid viscosity increases continuously through a large range in the transformation to an amorphous solid. Viscosity, expressed in units of poise, is used in Table 4 to specify characteristic working temperatures in the processing of the liquid precursors

of various oxide glasses. A poise is the centimetre-gram-second (cgs) unit of viscosity. It expresses the force needed to maintain a unit velocity difference between parallel plates separated by one centimetre of fluid: one poise equals one dyne-second per square centimetre. Molten glass may have a viscosity of 10¹³ poise (similar to honey on a cold day), and it quickly gets stiffer when cooled since the viscosity steeply increases with decreasing temperature. The ability to "tune" the viscosity of the melt (by changing temperature) allows glass to be conveniently processed and worked into desired shapes; glassblowing is a classic example of the usefulness of this widely exploited property.

Table 5 lists some important technological uses of amorphous solids. In addition to the application, the general type of amorphous solid used, and the material properties that make the application possible, the table also includes information about the chemical compositions of typical materials employed in these techniques. While the first entry—namely, window glass—represents the present status of a centuries-old technology, the other entries correspond to technologies that have blossomed during the second half of the 20th century. A significant theme of Table 5 is the role of amorphous solids in applications calling for large-area sheets or films. Amorphous solids often have great advantages over crystalline solids in such applications, since their use avoids the functional problems associated with polycrystallinity or the expense of preparing large single crystals. Thus, while it would be prohibitively expensive to fabricate large windows out of crystalline SiO₂ (quartz), it is practical to do so using SiO₂-based silicate glasses.

Transparent glasses. The terms glass and window glass are often used interchangeably in everyday language, so

Viscosity of glasses

Table 5: Some Technological Applications of Amorphous Solids

type of amorphous solid	representative material	application	special properties
Oxide glass	(SiO ₂) _{0.8} (Na ₂ O) _{0.2}	window glass	transparency, solidity, formability as large sheets
Oxide glass	(SiO ₂) _{0.9} (GeO ₂) _{0.1}	fibre-optic waveguides for communications networks	ultratransparency, purity, formability as uniform fibres
Organic polymer	polystyrene	structural materials, plastics	strength, light weight, ease of processing
Chalcogenide glass	Se, As ₂ Se ₃	copiers and laser printers	photoconductivity, formability as large-area films
Amorphous semiconductor	Si _{0.9} H _{0.1}	solar cells, copiers, flat-panel displays	photovoltaic optical properties, large-area thin films, semiconducting properties
Metallic glass	Fe _{0.8} B _{0.2}	transformer cores	ferromagnetism, low power loss, formability as long ribbons

familiar is this ancient architectural application of amorphous solids. Not only are oxide glasses, such as those characterized in Table 4, excellent for letting light in, they are also good for keeping cold out, because (as mentioned above) they are efficient thermal insulators.

The second application in Table 5 represents a modern development that carries the property of optical transparency to a phenomenal level. The transparency of the extraordinarily pure glasses that have been developed for fibre-optic telecommunications is so great that, at certain wavelengths, light can pass through 1 kilometre (0.6 mile) of glass and still retain 95 percent of its original intensity.

Glass fibres (transmitting optical signals) are now doing what copper wires (transmitting electrical signals) once did and are doing it more efficiently: carrying telephone messages around the planet. How this is done is schematically indicated in Figure 36. Digital electrical pulses produced by encoding of the voice-driven electrical signal are converted into light pulses by a semiconductor laser coupled to one end of the optical fibre. The signal is then transmitted over a long length of fibre as a stream of light pulses. At the far end it is converted back into electrical pulses and then into sound.

The glass fibre is somewhat thinner than a human hair. The simplest type, as sketched in the upper left of the figure, has a central core of ultratransparent glass surrounded by a coaxial cladding of a glass having a lower refractive index, n . This ensures that light rays propagating within the core, at small angles relative to the fibre axis, do not leak out but instead are 100 percent reflected at the core-cladding interface by the optical effect known as total internal reflection.

The great advantage provided by the substitution of light-transmitting fibres of ultratransparent oxide glass for electricity-transmitting wires of crystalline copper is that a single optical fibre can carry many more simultaneous conversations than can a thick cable packed with copper wires. This is the case because light waves oscillate at enormously high frequencies (about 2×10^{14} cycles per second for the infrared light generally used for fibre-optic telecommunications). This allows the light-wave signal carrier to be modulated at very high frequencies and to transmit a high volume of information traffic. Fibre-optic telecommunications have greatly expanded the information-transmitting capacity of the world's telecommunications networks.

Polymeric structural materials. Polystyrene, the organic polymer listed in Table 5, is a prototypical example of a polymeric glass. These glasses, whose atomic-scale structure has been discussed in connection with Figure 35B, make up a broad class of lightweight structural materials important in the automotive, aerospace, and construction industries. These materials are also ubiquitous in everyday experience as plastic molded objects. The quantity of polymer materials produced each year, measured in terms of volume, exceeds the quantity of steel produced.

Polystyrene is among the most important of the thermo-

plastic materials that, when heated (to the vicinity of the glass transition temperature), soften and flow controllably, enabling them to be processed at high speeds and on a large scale in the manufacture of molded products. The chemical formula of a polystyrene chain may be written as $(\text{CH}_2\text{CHC}_6\text{H}_5)_N$. The building block (inside the parentheses) consists of two backbone carbon atoms to which three hydrogen atoms and one phenyl (C_6H_5) ring are bonded as side groups. The polymerization index N reaches values above 10^5 . Polystyrene is a purely hydrocarbon polymer (*i.e.*, it contains only hydrogen and carbon); most organic polymers contain additional chemical components.

Amorphous semiconductors in electronics. Amorphous semiconductors, in the form of thin films prepared by methods such as that shown in Figure 32D, are important in applications requiring large areas of electronically active material. The first electronic application of amorphous semiconductors to occur on a large scale was in xerography (or electrostatic imaging), the process that provides the basis of plain-paper copiers. Amorphous selenium (Se) and, later, amorphous arsenic selenide (As_2Se_3) were used to form the thin-film, large-area photoconducting element that lies at the heart of the xerographic process. The photoconductor, which is an electrical insulator in the absence of light but which conducts electricity when illuminated, is exposed to an image of the document to be copied. Throughout the world—in offices, libraries, schools, and so forth—the xerographic process makes more than five billion copies every day. This process is also widely used in laser printers, in which the photoconductor is exposed to a digitally controlled on-and-off laser beam that is raster scanned (like the electron beam in a television tube) over the photoconductor surface.

Although still in use, selenium and arsenic selenide have been joined by other amorphous materials in this important technology. Polymeric organic glasses, in the form of thin films, are now used in multilayer photoconductor configurations in which the light is absorbed in one layer and electrical charge is transported through an adjacent layer. Both layers are formed of amorphous polymer films, and these photoreceptors can be made in the form of flexible belts.

Amorphous silicon thin films are used in solar cells that power handheld calculators. This important amorphous semiconductor is also used as the image sensor in facsimile ("fax") machines, and it serves as the photoreceptor in some xerographic copiers. All these applications exploit the ability of amorphous silicon to be vapour-deposited in the form of large-area thin films. As shown in Table 5, the practical form of this amorphous semiconductor is not pure silicon but a silicon-hydrogen alloy containing 10 percent hydrogen. The key role played by hydrogen, in what is now called hydrogenated amorphous silicon, emerged in a scientific puzzle that took years to solve. Stated briefly, hydrogen eliminates the electronic defects that are intrinsic to pure amorphous silicon.

Hydrogenated amorphous silicon also is used in high-

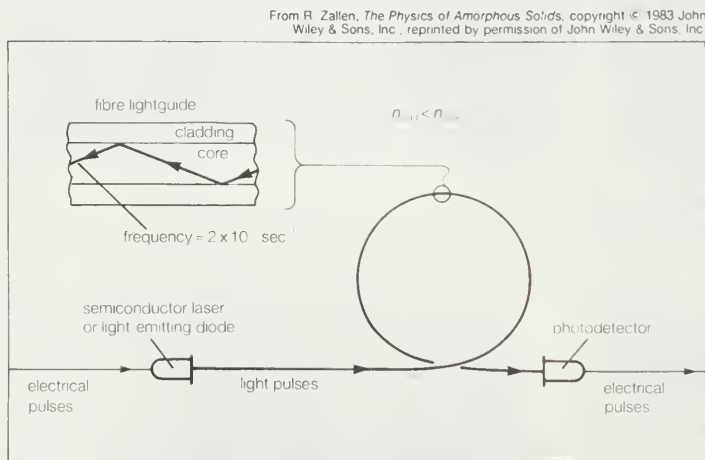


Figure 36: The use of ultratransparent glass fibres in telecommunications networks.

resolution flat-panel displays for computer monitors and for television screens. In such applications the large-area amorphous-semiconductor thin film is etched into an array of many tiny units, each of which forms the active element of a transistor that electronically turns on or off a small pixel (picture element) of a liquid-crystal display.

Magnetic glasses. The last entry in Table 5 is an application of metallic glasses having magnetic properties. These are typically iron-rich amorphous solids with compositions such as $\text{Fe}_{0.8}\text{B}_{0.2}$ iron-boron and $\text{Fe}_{0.8}\text{B}_{0.1}\text{Si}_{0.1}$ iron-boron-silicon. They are readily formed as long metallic glass ribbons by melt spinning or as wide sheets by planar flow casting. Ferromagnetic glasses are mechanically hard materials, but they are magnetically soft, meaning that they are easily magnetized by small magnetic fields. Also, because of their disordered atomic-scale structure,

they have higher electrical resistance than conventional (crystalline) magnetic materials. The three attributes of ease of manufacture, magnetic softness, and high electrical resistance make magnetic glasses extremely suitable for use in the magnetic cores of electrical power transformers. High electrical resistance (which arises here as a direct consequence of amorphicity) is a crucial property in this application, because it minimizes unwanted electrical eddy currents and cuts down on power losses. For these reasons, sheets of iron-based magnetic glasses are used as transformer-core laminations in electrical power applications.

Thin films of magnetic glass are finding use in many other applications. These include magnetic recording media for audio and video digital recording, as well as recording heads used with magnetic disks. (R.Z.)

Electrical power applications

LIQUID STATE

The most obvious physical properties of a liquid are its retention of volume and its conformation to the shape of its container. When a liquid substance is poured into a vessel, it takes the shape of the vessel, and, as long as the substance stays in the liquid state, it will remain inside the vessel. Furthermore, when a liquid is poured from one vessel to another, it retains its volume (as long as there is no vaporization or change in temperature) but not its shape. These properties serve as convenient criteria for distinguishing the liquid state from the solid and gaseous states. Gases, for example, expand to fill their container so that the volume they occupy is the same as that of the container. Solids retain both their shape and volume when moved from one container to another.

Liquids may be divided into two general categories: pure liquids and liquid mixtures. On Earth, water is the most abundant liquid, although much of the water with which organisms come into contact is not in pure form but is a mixture in which various substances are dissolved. Such mixtures include those fluids essential to life—blood, for example—beverages, and seawater. Seawater is a liquid mixture in which a variety of salts have been dissolved in water. Even though in pure form these salts are solids, in oceans they are part of the liquid phase. Thus, liquid mixtures contain substances that in their pure form may themselves be liquids, solids, or even gases.

The liquid state sometimes is described simply as the state that occurs between the solid and gaseous states, and for simple molecules this distinction is unambiguous. However, clear distinction between the liquid, gaseous, and solid states holds only for those substances whose molecules are composed of a small number of atoms. When the number exceeds about 20, the liquid may often be cooled below the true melting point to form a glass (see *Solid state: Amorphous solids*), which has many of the mechanical properties of a solid but lacks crystalline order. If the number of atoms in the molecule exceeds about 100–200, the classification into solid, liquid, and gas ceases to be useful. At low temperatures such substances are usually glasses or amorphous solids, and their rigidity falls with increasing temperature—*i.e.*, they do not have fixed melting points; some may, however, form true liquids. With these large molecules, the gaseous state is not attainable, because they decompose chemically before the temperature is high enough for the liquid to evaporate. Synthetic and natural high polymers (*e.g.*, nylon and rubber) behave in this way.

If the molecules are large, rigid, and either roughly planar or linear, as in cholesteryl acetate or *p*-azoxyanisole, the solid may melt to an anisotropic liquid (*i.e.*, one that is not uniform in all directions) in which the molecules are free to move about but have great difficulty in rotating. Such a state is called a liquid crystal, and the anisotropy produces changes of the refractive index (a measure of the change in direction of light when it passes from one medium into another) with the direction of the incident light and hence leads to unusual optical effects. Liquid crystals have found widespread applications in temper-

ature-sensing devices and in displays for watches and calculators. However, no inorganic compounds and only about 5 percent of the known organic compounds form liquid crystals. The theory of normal liquids is, therefore, predominantly the theory of the behaviour of substances consisting of simple molecules.

A liquid lacks both the strong spatial order of a solid, though it has the high density of solids, and the absence of order of a gas that results from the low density of gases—*i.e.*, gas molecules are relatively free of each other's influence. The combination of high density and of partial order in liquids has led to difficulties in developing quantitatively acceptable theories of liquids. Understanding of the liquid state, as of all states of matter, came with the kinetic molecular theory, which stated that matter consisted of particles in constant motion and that this motion was the manifestation of thermal energy. The greater the thermal energy of the particle, the faster it moved.

In very general terms, the particles that constitute matter include molecules, atoms, ions, and electrons. In a gas these particles are far enough from one another and are moving fast enough to escape each other's influence, which may be of various kinds—such as attraction or repulsion due to electrical charges and specific forces of attraction that involve the electrons orbiting around atomic nuclei. The motion of particles is in a straight line, and the collisions that result occur with no loss of energy, although an exchange of energies may result between colliding particles. When a gas is cooled, its particles move more slowly, and those slow enough to linger in each other's vicinity will coalesce, because a force of attraction will overcome their lowered kinetic energy and, by definition, thermal energy. Each particle, when it joins others in the liquid state, gives up a measure of heat called the latent heat of liquefaction, but each continues to move at the same speed within the liquid as long as the temperature remains at the condensation point. The distances that the particles can travel in a liquid without colliding are on the order of molecular diameters. As the liquid is cooled, the particles move more slowly still, until at the freezing temperature the attractive energy produces so high a density that the liquid freezes into the solid state. They continue to vibrate, however, at the same speed as long as the temperature remains at the freezing point, and their latent heat of fusion is released in the freezing process. Heating a solid provides the particles with the heat of fusion necessary to allow them to escape one another's influence enough to move about in the liquid state. Further heating provides the liquid particles with their heat of evaporation, which enables them to escape one another completely and enter the vapour, or gaseous, state.

This starkly simplified view of the states of matter ignores many complicating factors, the most important being the fact that no two particles need be moving at the same speed in a gas, liquid, or solid and the related fact that even in a solid some particles may have acquired the energy necessary to exist as gas particles, while even in a gas some particles may be practically motionless for a

Transition to the liquid state

brief time. It is the average kinetic energy of the particles that must be considered, together with the fact that the motion is random. At the interface between liquid and gas and between liquid and solid, an exchange of particles is always taking place: slow gas molecules condensing at the liquid surface and fast liquid molecules escaping into the gas. An equilibrium state is reached in any closed system, so that the number of exchanges in either direction is the same. Because the kinetic energy of particles in the liquid state can be defined only in statistical terms (*i.e.*, every possible value can be found), discussion of the liquid (as well as the gaseous) state at the molecular level involves formulations in terms of probability functions.

Behaviour of pure liquids

When the temperature and pressure of a pure substance are fixed, the equilibrium state of the substance is also fixed. This is illustrated in Figure 37, which shows the phase diagram for pure argon. In the diagram a single phase is shown as an area, two as a line, and three as the intersection of the lines at the triple point, *T*. Along the line *TC*, called the vapour-pressure curve, liquid and vapour exist in equilibrium. The liquid region exists to the left and above this line while the gas, or vapour, region exists below it. At the upper extreme, this curve ends at the critical point, *C*. If line *TC* is crossed by moving directly from point *P* to *S*, there is a distinct phase change accompanied by abrupt changes in the physical properties of the substance (*e.g.*, density, heat capacity, viscosity, and dielectric constant) because the vapour and liquid phases have distinctly different properties. At the critical point, however, the vapour and liquid phases become identical, and above the critical point, the two phases are no longer distinct. Thus, if the substance moves from point *P* to *S* by the path *PQRS* so that no phase-change lines are crossed, the change in properties will be smooth and continuous, and the specific moment when the substance converts from a liquid to a gas is not clearly defined. In fact, the path *PQRS* demonstrates the essential continuity of state between liquid and gas, which differ in degree but which together constitute the single fluid state. Strictly speaking, the term liquid should be applied only to the denser of the two phases on the line *TC*, but it is generally extended to any dense fluid state at low temperatures—*i.e.*, to the area lying within the angle *CTM*.

The extension of line *TC* below the triple point is called the sublimation curve. It represents the equilibrium between solid and gas, and when the sublimation curve is crossed, the substance changes directly from solid to gas. This conversion occurs when dry ice (solid carbon dioxide) vaporizes at atmospheric pressure to form gaseous carbon dioxide because the triple-point pressure for carbon dioxide is greater than atmospheric pressure. Line *TM* is the melting curve and represents an equilibrium between solid and liquid; when this curve is crossed from left to right, solid changes to liquid with the associated abrupt change in properties.

The melting curve is initially much steeper than the vapour-pressure curve; hence, as the pressure is changed, the temperature does not change much, and the melting temperature is little affected by pressure. No substance has been found to have a critical point on this line, and there are theoretical reasons for supposing that it continues indefinitely to high temperatures and pressures, until the substance is so compressed that the molecules break up into atoms, ions, and electrons. At pressures above 10^6 bars (one bar is equal to 0.987 atmosphere, where one atmosphere is the pressure exerted by the air at sea level), it is believed that most substances pass into a metallic state.

It is possible to cool a gas at constant pressure to a temperature lower than that of the vapour-pressure line without producing immediate condensation, since the liquid phase forms readily only in the presence of suitable nuclei (*e.g.*, dust particles or ions) about which the drops can grow. Unless the gas is scrupulously cleaned, such nuclei remain; a subcooled vapour is unstable and will ultimately condense. It is similarly possible to superheat a liquid to a temperature where, though still a liquid, the gas

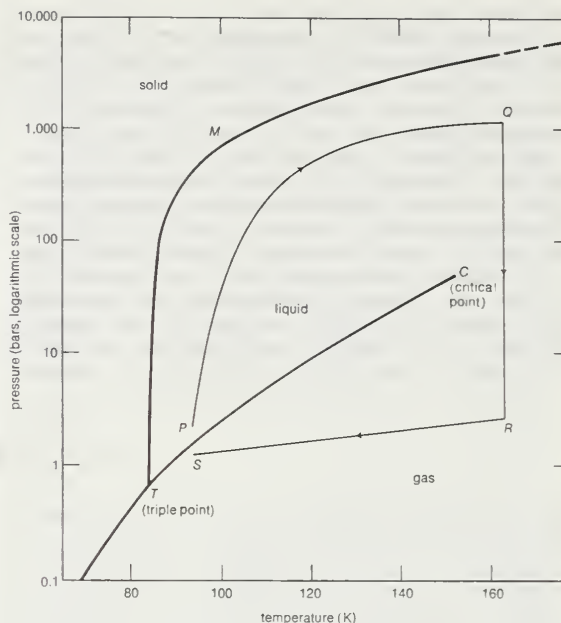


Figure 37: Phase diagram of argon. Heavy lines show the boundaries between liquid and gas (*TC*) and between liquid and solid (*TM*). Curve *PQRS* describes a path from a typically liquid state (*P*) to a typically gaseous state (*S*; see text).

is the stable phase. Again, this occurs most readily with clean liquids heated in smooth vessels, because bubble formation occurs around foreign particles or sharp points. When the superheated liquid changes to gas, it does so with almost explosive violence. A liquid also may be subcooled to below its freezing temperature.

To a certain extent the behaviour of all substances is similar to that described in Figure 37. The parameters that vary from substance to substance are the particular values of the triple-point and critical-point temperature and pressure, the size of the various regions, and the slopes of the lines. Triple-point temperatures range from 14 K (0 K equals -273.15°C [-459.67°F]), for hydrogen to temperatures too high for accurate measurement. Triple-point pressures are generally low, that of carbon dioxide at 5.2 bars being one of the highest. Most are around 10^{-3} bar, and those of some hydrocarbons are as low as 10^{-7} bar. The normal melting point of a substance is defined as the melting temperature at a pressure of one atmosphere (equivalent to 1.01325 bars); it differs little from the triple-point temperature, because of the steepness of melting lines (*TM* in Figure 37). Critical temperatures (the maximum temperature at which a gas can be liquefied by pressure) range from 5.2 K, for helium, to temperatures too high to measure. Critical pressures (the vapour pressure at the critical temperature) are generally about 40–100 bars. The normal boiling point is the temperature at which the vapour pressure reaches one atmosphere. The normal liquid range is defined as the temperature interval between the normal melting point and the normal boiling point, but such a restriction is artificial, the true liquid range being from triple point to critical point. Substances whose triple-point pressures are above atmospheric (*e.g.*, carbon dioxide) have no normal liquid range but sublime at atmospheric pressure.

Each of the three two-phase lines in Figure 37 can be described by the Clapeyron equation:

$$\frac{dp}{dT} = \frac{\Delta H}{T\Delta V} \quad (1)$$

In this equation, dp/dT is the slope of the curve under consideration—*i.e.*, either the melting, sublimation, or vapour-pressure curve. ΔH is the latent heat required for the phase change, and ΔV is the change in volume associated with the phase change. Thus, for the sublimation and vapour-pressure curves, since ΔH and ΔV are both positive (*i.e.*, heat is required for vaporization, and the volume increases on vaporization), the slope is always

Triple-point temperature

positive. Although not apparent from Figure 37, the slope of the sublimation curve immediately below the triple point is greater than the slope of the vapour-pressure curve immediately above it, so that the vapour-pressure curve is not continuous through the triple point. This is consistent with equation (1) because the heat of sublimation for a substance is somewhat larger than its heat of vaporization. The slope of the melting line is usually positive, but there are a few substances, such as water and bismuth, for which the melting-line slope is negative. The negative melting-line slope is consistent with equation (1) because, for these two substances, the density of the solid is less than the density of the liquid. This is the reason ice floats. For water, this negative volume change (*i.e.*, shrinking) persists to 2.1 kilobars and -22°C , at which point the normal form of ice changes to a denser form, and thereafter the change in volume on melting is positive.

At the critical point the liquid is identical to the vapour phase, and near the critical point the liquid behaviour is somewhat similar to vapour-phase behaviour. While the particular values of the critical temperature and pressure vary from substance to substance, the nature of the behaviour in the vicinity of the critical point is similar for all compounds. This fact has led to a method that is commonly referred to as the law of corresponding states. Roughly speaking, this approach presumes that, if the phase diagram is plotted using reduced variables, the behaviour of all substances will be more or less the same. Reduced variables are defined by dividing the actual variable by its associated critical constant: the reduced temperature, T_r , equals T/T_c , and the reduced pressure, p_r , equals p/p_c . Then for all substances the critical point occurs at a value of T_r and p_r equal to unity. This approach has been used successfully to develop equations to correlate and predict a number of liquid-phase properties including vapour pressures, saturated and compressed liquid densities, heat capacities, and latent heats of vaporization. The corresponding states approach works remarkably well at temperatures between the normal boiling point and the critical point for many compounds but tends to break down near and below the triple-point temperature. At these temperatures the liquid is influenced more by the behaviour of the solid, which has not been successfully correlated by corresponding states methods.

Many of the properties of a liquid near its triple point are closer to those of the solid than to those of the gas. It has a high density (typically 0.5–1.5 grams per cubic centimetre [0.02–0.05 pound per cubic inch]), a high refractive index (which varies from 1.3 to 1.8 for liquids), a high heat capacity at constant pressure (two to four joules per gram per kelvin, one joule being equal to 0.239 calorie), and a low compressibility ($0.5\text{--}1 \times 10^{-4}$ per bar). The compressibility falls to values characteristic of a solid (0.1×10^{-4} per bar or less) as the pressure increases. A simple and widely used equation describes the change of specific volume with pressure. If $V(p)$ is the volume at pressure p , $V(0)$ is volume at zero pressure, and A and B are positive parameters (constants whose values may be arbitrarily assigned), then the difference in volume resulting from a change in pressure equals the product of A , the pressure, and the volume at zero pressure, divided by the sum of B and the pressure. This is written:

$$V(0) - V(p) = \frac{ApV(0)}{(B + p)}. \quad (2)$$

The pressure parameter B is close to the pressure at which the compressibility has fallen to half its initial value and is generally about 500 bars for liquids near their triple points. It falls rapidly with increasing temperature.

As a liquid is heated along its vapour-pressure curve, TC , its density falls and its compressibility rises. Conversely, the density of the saturated vapour in equilibrium with the liquid rises; *i.e.*, the number of gas molecules in a fixed space above the liquid increases. Liquid and gas states approach each other with increasing rapidity as the temperature approaches C , until at this point they become identical and have a density about one-third that of the liquid at point T . The change of saturated-gas density (ρ_g) and liquid density (ρ_l) with temperature T can be ex-

pressed by a simple equation when the temperature is close to critical. If ρ_c is the density at the critical temperature T_c , then the difference between densities equals the difference between temperatures raised to a factor called beta, β :

$$(\rho_l - \rho_c) = (\rho_c - \rho_g) = (T_c - T)^\beta, \quad (3)$$

where β is about 0.34. The compressibility and the heat capacity of the gas at constant pressure (C_p) become infinite as T approaches T_c from above along the path of constant density. The infinite compressibility implies that the pressure no longer restrains local fluctuations of density. The fluctuations grow to such an extent that their size is comparable with the wavelength of light, which is therefore strongly scattered. Hence, at the critical point, a normally transparent liquid is almost opaque and usually dark brown in colour. The classical description of the critical point and the results of modern measurement do not agree in detail, but recent considerations of thermodynamic stability show that there are certain regularities in behaviour that are common to all substances.

Between a liquid and its corresponding vapour there is a dividing surface that has a measurable tension; work must be done to increase the area of the surface at constant temperature. Hence, in the absence of gravity or during free fall, the equilibrium shape of a volume of liquid is one that has a minimum area—*i.e.*, a sphere. In the Earth's field this shape is found only for small drops, for which the gravitational forces, since they are proportional to the volume, are negligible compared with surface forces, which are proportional to the area. The surface tension falls with rising temperature and vanishes at the critical point. There is a similar dividing surface between two immiscible liquids, but this usually has lower tension. It is believed that there is a tension also between a liquid and a solid, though it is not directly measurable because of the rigidity of the solid; it may be inferred, however, under certain assumptions, from the angle of contact between the liquid and the solid (*i.e.*, the angle at which the liquid's surface meets the solid). If this angle is zero, the liquid surface is parallel to the solid surface and is said to wet the solid completely.

MOLECULAR STRUCTURE OF LIQUIDS

For a complete understanding of the liquid state of matter, an understanding of behaviour on the molecular level is necessary. Such behaviour is characterized by two quantities called the intermolecular pair potential function, u , and the radial distribution function, g . The pair potential gives information about the energy due to the interaction of a pair of molecules and is a function of the distance r between their centres. Information about the structure or the distances between pairs of molecules is contained in the radial distribution function. If g and u are known for a substance, macroscopic properties can be calculated.

In an ideal gas—where there are no forces between molecules, and the volume of the molecules is negligible— g is unity, which means that the chance of encountering a second molecule when moving away from a central molecule is independent of position. In a solid, g takes on discrete values at distances that correspond to the locations associated with the solid's crystal structure. Liquids possess neither the completely ordered structure of a solid crystal nor the complete randomness of an ideal gas. The structure in a liquid is intermediate to these two extremes—*i.e.*, the molecules in a liquid are free to move about, but there is some order because they remain relatively close to one another. Although there are an infinite number of possible positions one molecule may assume with respect to another, some are more likely than others. This is illustrated by Figure 38, which shows an example of the radial distribution function for the dense packing typical of liquids. In this figure, g is a measure of the probability of finding the centre of one molecule at a distance r from the centre of a second molecule. For values of r less than those of the molecular diameter, d , g goes to zero. This is consistent with the fact that two molecules cannot occupy the same space. The most likely location for a second molecule with respect to a central molecule is slightly more than one molecular diameter away, which

Law of corresponding states

Shape of drops

Radial distribution function

reflects the fact that in liquids the molecules are packed almost against one another. The second most likely location is a little more than two molecular diameters away, but beyond the third layer preferred locations damp out, and the chance of finding the centre of a molecule becomes independent of position.

The pair potential function, u , is a large positive number for r less than d , assumes a minimum value at the most preferred location (this corresponds to the maximum of the curve in Figure 38), and damps out to zero as r approaches infinity. The large positive value of u corresponds to a strong repulsion, while the minimum corresponds to the net result of repulsive and attractive forces.

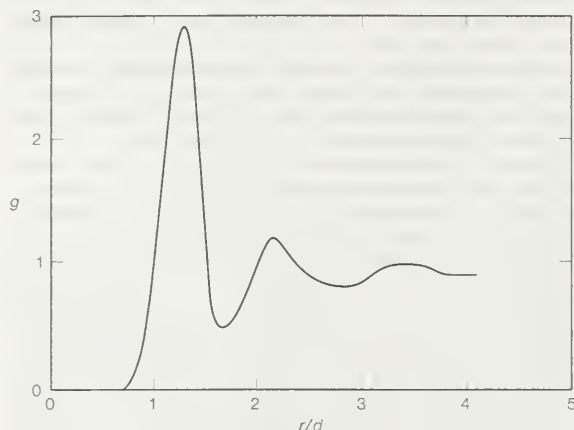


Figure 38: Radial distribution function for a dense fluid. The radial distribution function, g , is a measure of the probability of finding the centre of one molecule at a distance r from a central molecule. The molecular diameter is d .

There are two methods of measuring the radial distribution function g : first, by X-ray or neutron diffraction from simple fluids and, second, by computer simulation of the molecular structure and motions in a liquid. In the first, the liquid is exposed to a specific, single wavelength (monochromatic) radiation, and the observed results are then subjected to a mathematical treatment known as a Fourier transform.

The second method of obtaining the radial distribution function g supposes that the energy of interaction, u , for the liquid under study is known. A computer model of a liquid is set up, in which between 100 and 1,000 molecules are contained within a cube. There are now two methods of proceeding: by Monte Carlo calculation or by what is called molecular dynamics; only the latter is discussed here. Each molecule is assigned a random position and velocity, and Newton's equations of motion are solved to calculate the path of each molecule in the changing field of all the others. A molecule that leaves the cell is deemed to be replaced by a new one with equal velocity entering through the corresponding point on the opposite wall. After a few collisions per molecule, the distribution of velocities conforms with equations worked out by the Scottish physicist James Clerk Maxwell, and after a longer time the mean positions are those appropriate to the density and mean kinetic energy (*i.e.*, temperature) of the liquid. Functions such as the radial distribution function g can now be evaluated by taking suitable averages as the system evolves in time. Since 1958 such computer experiments have added more to the knowledge of the molecular structure of simple liquids than all the theoretical work of the previous century and continue to be an active area of research for not only pure liquids but liquid mixtures as well.

SPEED OF SOUND AND ELECTRIC PROPERTIES

A sound wave is a series of longitudinal compressions and expansions that travels through a liquid at a speed of about one kilometre per second (0.62 mile per second), or about three times the speed of sound in air. If the frequency is not too high, the compressions and expansions are adiabatic (*i.e.*, the changes take place without transfer of heat) and reversible. Conduction of energy from the hot

(compressed) to the cold (expanded) regions of the liquid introduces irreversible effects, which are dissipative, and thus such conduction leads to the absorption of the sound. A longitudinal compression (in the direction of the wave) is a combination of a uniform compression and a shearing stress (a force that causes one plane of a substance to glide past an adjacent plane). Hence, both bulk and shear viscosity also govern the propagation of sound in a liquid.

If a liquid is placed in a static electric field, the field exerts a force on any free carriers of electric charge in the liquid, and the liquid, therefore, conducts electricity. Such carriers are of two kinds: mobile electrons and ions. The former are present in abundance in liquid metals, which have conductivities that are generally about one-third of the conductivity of the corresponding solid. The decrease in conductivity upon melting arises from the greater disorder of the positive ions in the liquid and hence their greater ability to scatter electrons. The contribution of the ions is small, less than 5 percent in most liquid metals, but it is the sole cause of conductivity in molten salts and in their aqueous solutions. Such conductivities vary widely but are much lower than those of liquid metals.

Nonionic liquids (those composed of molecules that do not dissociate into ions) have negligible conductivities, but they are polarized by an electric field: that is, the liquid develops positive and negative poles and also a dipole moment (which is the product of the pole strength and the distance between the poles) that is oriented against the field, from which the liquid acquires energy. This polarization is of three kinds: electron, atomic, and orientation. In electron polarization the electrons in each atom are displaced from their usual positions, giving each molecule a small dipole moment. The contribution of electron polarization to the dielectric constant (see below *Solutions and solubilities: Classes of solutions: Electrolytes and nonelectrolytes*) of the liquid is numerically equal to the square root of its refractive index. The second effect, atomic polarization, arises because there is a relative change in the mean positions of the atomic nuclei within the molecules. This generally small effect is observed at radio frequencies but not at optical, and so it is missing from the refractive index. The third effect, orientation polarization, occurs with molecules that have permanent dipole moments. These molecules are partially aligned by the field and contribute heavily to the polarization. Thus, the dielectric constant of a nonpolar liquid, such as a hydrocarbon, is about 2, that of a weakly polar liquid, such as chloroform or ethyl ether, about 5, while those of highly polar liquids, such as ethanol and water, range from 25 to 80.

(Jo.S.R./B.E.P.)

Solutions and solubilities

The ability of liquids to dissolve solids, other liquids, or gases has long been recognized as one of the fundamental phenomena of nature encountered in daily life. The practical importance of solutions and the need to understand their properties have challenged numerous writers since the Ionian philosophers and Aristotle. Though many physicists and chemists have devoted themselves to a study of solutions, as of the early 1990s it was still an incompletely understood subject under active investigation.

A solution is a mixture of two or more chemically distinct substances that is said to be homogeneous on the molecular scale—the composition at any one point in the mixture is the same as that at any other point. This is in contrast to a suspension (or slurry), in which small discontinuous particles are surrounded by a continuous fluid. Although the word solution is commonly applied to the liquid state of matter, solutions of solids and gases are also possible; brass, for example, is a solution of copper and zinc, and air is a solution primarily of oxygen and nitrogen with a few other gases present in relatively small amounts.

The ability of one substance to dissolve another depends always on the chemical nature of the substances, frequently on the temperature, and occasionally on the pressure. Water, for example, readily dissolves methyl alcohol but does not dissolve mercury; it barely dissolves benzene at room temperature but does so increasingly as the temperature

Nonionic liquids

Atomic polarization

Use of a computer model

rises. While the solubility in water of the gases present in air is extremely small at atmospheric pressure, it becomes appreciable at high pressures where, in many cases, the solubility of a gas is (approximately) proportional to its pressure. Thus, a diver breathes air (four-fifths nitrogen) at a pressure corresponding to the pressure around him, and, as he goes deeper, more air dissolves in his blood. If he ascends rapidly, the solubility of the gases decreases so that they leave his blood suddenly, forming bubbles in the blood vessels. This condition (known as the bends) is extremely painful and may cause death; it can be alleviated by breathing, instead of air, a mixture of helium and oxygen because the solubility of helium in blood is much lower than that of nitrogen.

The solubility of one fluid in another may be complete or partial; thus, at room temperature water and methyl alcohol mix in all proportions, but 100 grams of water dissolve only 0.07 gram of benzene. Though it is generally supposed that all gases are completely miscible—*i.e.*, mutually soluble in all proportions—this is true only at normal pressures. At high pressures pairs of chemically unlike gases may exhibit only limited miscibility; for example, at 20° C helium and xenon are completely miscible at pressures below 200 atmospheres but become increasingly immiscible as the pressure rises.

The ability of a liquid to dissolve selectively forms the basis of common separation operations in chemical and related industries. A mixture of two gases, carbon dioxide and nitrogen, can be separated by bringing it into contact with ethanalamine, a liquid solvent that readily dissolves carbon dioxide but barely dissolves nitrogen. In this process, called absorption, the dissolved carbon dioxide is later recovered, and the solvent is made usable again by heating the carbon dioxide-rich solvent, since the solubility of a gas in a liquid usually (but not always) decreases with rising temperature. A similar absorption operation can remove a pollutant such as sulfur dioxide from smokestack gases in a plant using sulfur-containing coal or petroleum as fuel.

The process wherein a dissolved substance is transferred from one liquid to another is called extraction. As an example, phenolic pollutants (organic compounds of the types known as phenol, cresol, and resorcinol) are frequently found in industrial aqueous waste streams, and, since these phenolics are damaging to marine life, it is important to remove them before sending the waste stream back to a lake or river. One such removal technique is to bring the waste stream into contact with a water-insoluble solvent (*e.g.*, an organic liquid such as a high-boiling hydrocarbon) that has a strong affinity for the phenolic pollutant. The solubility of the phenolic in the solvent divided by that in water is called the distribution coefficient, and it is clear that for an efficient extraction process it is desirable to have as large a distribution coefficient as possible.

CLASSES OF SOLUTIONS

Electrolytes and nonelectrolytes. Broadly speaking, liquid mixtures can be classified as either solutions of electrolytes or solutions of nonelectrolytes. Electrolytes are substances that can dissociate into electrically charged particles called ions, while nonelectrolytes consist of molecules that bear no net electric charge. Thus, when ordinary salt (sodium chloride, formula NaCl) is dissolved in water, it forms an electrolytic solution, dissociating into positive sodium ions (Na⁺) and negative chloride ions (Cl⁻), whereas sugar dissolved in water maintains its molecular integrity and does not dissociate. Because of its omnipresence, water is the most common solvent for electrolytes; the ocean is a solution of electrolytes. Electrolyte solutions, however, are also formed by other solvents (such as ammonia and sulfur dioxide) that have a large dielectric constant (a measure of the ability of a fluid to decrease the forces of attraction and repulsion between charged particles). The energy required to separate an ion pair (*i.e.*, one ion of positive charge and one ion of negative charge) varies inversely with the dielectric constant, and, therefore, appreciable dissociation into separate ions occurs only in solvents with large dielectric constants.

Most electrolytes (for example, salts) are nonvolatile, which means that they have essentially no tendency to enter the vapour phase. There are, however, some notable exceptions, such as hydrogen chloride (HCl), which is readily soluble in water, where it forms hydrogen ions (H⁺) and chloride ions (Cl⁻). At normal temperature and pressure, pure hydrogen chloride is a gas, and, in the absence of water or some other ionizing solvent, hydrogen chloride exists in molecular, rather than ionic, form.

Solutions of electrolytes readily conduct electricity, whereas nonelectrolyte solutions do not. A dilute solution of hydrogen chloride in water is a good electrical conductor, but a dilute solution of hydrogen chloride in a hydrocarbon is a good insulator. Because of the large difference in dielectric constants, hydrogen chloride is ionized in water but not in hydrocarbons.

Weak electrolytes. While classification under the heading electrolyte-solution or nonelectrolyte-solution is often useful, some solutions have properties near the boundary between these two broad classes. Although such substances as ordinary salt and hydrogen chloride are strong electrolytes—*i.e.*, they dissociate completely in an ionizing solvent—there are many substances, called weak electrolytes, that dissociate to only a small extent in ionizing solvents. For example, in aqueous solution, acetic acid can dissociate into a positive hydrogen ion and a negative acetate ion (CH₃COO⁻), but it does so to a limited extent; in an aqueous solution containing 50 grams acetic acid and 1,000 grams water, less than 1 percent of the acetic acid molecules are dissociated into ions. Therefore, a solution of acetic acid in water exhibits some properties associated with electrolyte solutions (*e.g.*, it is a fair conductor of electricity), but in general terms it is more properly classified as a nonelectrolyte solution. By similar reasoning, an aqueous solution of carbon dioxide is also considered a nonelectrolyte solution even though carbon dioxide and water have a slight tendency to form carbonic acid, which, in turn, dissociates to a small extent to hydrogen ions and bicarbonate ions (HCO₃⁻).

Endothermic and exothermic solutions. When two substances mix to form a solution, heat is either evolved (an exothermic process) or absorbed (an endothermic process); only in the special case of an ideal solution do substances mix without any heat effect. Most simple molecules mix with a small endothermic heat of solution, while exothermic heats of solution are observed when the components interact strongly with one another. An extreme example of an exothermic heat of mixing is provided by adding an aqueous solution of sodium hydroxide, a powerful base, to an aqueous solution of hydrogen chloride, a powerful acid; the hydroxide ions (OH⁻) of the base combine with the hydrogen ions of the acid to form water, a highly exothermic reaction that yields 75,300 calories per 100 grams of water formed. In nonelectrolyte solutions, heat effects are usually endothermic and much smaller, often about 100 calories, when roughly equal parts are mixed to form 100 grams of mixture.

Formation of a solution usually is accompanied by a small change in volume. If equal parts of benzene and stannic chloride are mixed, the temperature drops; if the mixture is then heated slightly to bring its temperature back to that of the unmixed liquids, the volume increases by about 2 percent. On the other hand, mixing roughly equal parts of acetone and chloroform produces a small decrease in volume, about 0.2 percent. It frequently happens that mixtures with endothermic heats of mixing expand—*i.e.*, show small increases in volume—while mixtures with exothermic heats of mixing tend to contract.

A large decrease in volume occurs when a gas is dissolved in a liquid. For example, at 0° C and atmospheric pressure, the volume of 28 grams of nitrogen gas is 22,400 cubic centimetres. When these 28 grams of nitrogen are dissolved in an excess of water, the volume of the water increases only 40 cubic centimetres; the decrease in volume accompanying the dissolution of 28 grams of nitrogen in water is therefore 22,360 cubic centimetres. In this case, it is said that the nitrogen gas has been condensed into a liquid, the word condense meaning "to make dense"—*i.e.*, to decrease the volume.

Volume changes caused by the formation of a solution

Effects of temperature and pressure

Ion pairs

PROPERTIES OF SOLUTIONS

Composition ratios. The composition of a liquid solution means the composition of that solution in the bulk—that is, of that part that is not near the surface. The interface between the liquid solution and some other phase (for example, a gas such as air) has a composition that differs, sometimes very much, from that of the bulk. The environment at an interface is significantly different from that throughout the bulk of the liquid, and in a solution the molecules of a particular component may prefer one environment over the other. If the molecules of one component in the solution prefer to be at the interface as opposed to the bulk, it is said that this component is positively adsorbed at the interface. In aqueous solutions of organic liquids, the organic component is usually positively adsorbed at the solution-air interface; as a result, it is often possible to separate a mixture of an organic solute from water by a process called froth separation. Air is bubbled vigorously into the solution, and a froth is formed. The composition of the froth differs from that of the bulk because the organic solute concentrates at the interfacial region. The froth is mechanically removed and collapsed, and, if further separation is desired, a new froth is generated. The tendency of some dissolved molecules to congregate at the surface has been utilized in water conservation. A certain type of alcohol, when added to water, concentrates at the surface to form a barrier to evaporating water molecules. In warm climates, therefore, water loss by evaporation from lakes can be significantly reduced by introducing a solute that adsorbs positively at the lake-air interface.

The composition of a solution can be expressed in a variety of ways, the simplest of which is the weight fraction, or weight percent; for example, the salt content of seawater is about 3.5 weight percent—*i.e.*, of 100 grams of seawater, 3.5 grams is salt. For a fundamental understanding of solution properties, however, it is often useful to express composition in terms of molecular units such as molecular concentration, molality, or mole fraction. To understand these terms, it is necessary to define atomic and molecular weights. The atomic weight of elements is a relative figure, with one atom of the carbon-12 isotope being assigned the atomic weight of 12; the atomic weight of hydrogen is then approximately 1, of oxygen approximately 16, and the molecular weight of water (H_2O) 18. The atomic and molecular theory of matter asserts that the atomic weight of any element in grams must contain the same number of atoms as the atomic weight in grams (the gram-atomic weight) of any other element. Thus, two grams of molecular hydrogen (H_2)—its gram-molecular weight—contain the same number of molecules as 18 grams of water or 32 grams of oxygen molecules (O_2). Further, a specified volume of any gas (at low pressure) contains the same number of molecules as the same volume of any other gas at the same temperature and pressure. At standard temperature and pressure ($0^\circ C$ and one atmosphere) the volume of one gram-molecular weight of any gas has been determined experimentally to be approximately 22.4 litres (23.7 quarts). The number of molecules in this volume of gas, or in the gram-molecular weight of any compound, is called Avogadro's number.

Molarity. Molecular concentration is the number of molecules of a particular component per unit volume. Since the number of molecules in a litre or even a cubic centimetre is enormous, it has become common practice to use what are called molar, rather than molecular, quantities. A mole is the gram-molecular weight of a substance and, therefore, also Avogadro's number of molecules (6.02×10^{23}). Thus, the number of moles in a sample is the weight of the sample divided by the molecular weight of the substance; it is also the number of molecules in the sample divided by Avogadro's number. Instead of using molecular concentration, it is more convenient to use molar concentration: instead of saying, for example, that the concentration is 12.04×10^{23} molecules per litre, it is simpler to say that it is two moles per litre. Concentration in moles per litre (*i.e.*, molarity) is usually designated by the letter M.

Molality. In electrolyte solutions it is common to dis-

tinguish between the solvent (usually water) and the dissolved substance, or solute, which dissociates into ions. For these solutions it is useful to express composition in terms of molality, designated as m , a unit proportional to the number of undissociated solute molecules (or, alternatively, to the number of ions) per 1,000 grams of solvent. The number of molecules or ions in 1,000 grams of solvent usually is very large, so molality is defined as the number of moles per 1,000 grams of solvent.

Formality. Many compounds do not exist in molecular form, either as pure substances or in their solutions. The particles that make up sodium chloride ($NaCl$), for example, are sodium ions (Na^+) and chloride ions (Cl^-), and, although equal numbers of these two ions are present in any sample of sodium chloride, no Na^+ ion is associated with a particular Cl^- ion to form a neutral molecule having the composition implied by the formula. Therefore, even though the compositions of such compounds are well defined, it would be erroneous to express concentrations of their solutions in terms of molecular weights. A useful concept in cases of this kind is that of the formula weight, defined as the sum of the weights of the atoms in the formula of the compound; thus, the formula weight of sodium chloride is the sum of the atomic weights of sodium and chlorine, 23 plus 35.5, or 58.5, and a solution containing 58.5 grams of sodium chloride per litre is said to have a concentration of one formal, or 1 F.

Mole fraction and mole percentage. It often is useful to express the composition of nonelectrolyte solutions in terms of mole fraction or mole percentage. In a binary mixture—*i.e.*, a mixture of two components, 1 and 2—there are two mole fractions, x_1 and x_2 , which satisfy the relation $x_1 + x_2 = 1$. The mole fraction x_1 is the fraction of molecules of species 1 in the solution, and x_2 is the fraction of molecules of species 2 in the solution. (Mole percentage is the mole fraction multiplied by 100.)

Volume fraction. The composition of a nonelectrolyte solution containing very large molecules, known as polymers, is most conveniently expressed by the volume fraction (Φ)—*i.e.*, the volume of polymer used to prepare the solution divided by the sum of that volume of polymer and the volume of the solvent.

Equilibrium properties. A quantitative description of liquid-solution properties when the system is in equilibrium is provided by relating the vapour pressure of the solution to its composition. The vapour pressure of a liquid, pure or mixed, is the pressure exerted by those molecules that escape from the liquid to form a separate vapour phase above the liquid. If a quantity of liquid is placed in an evacuated, closed container the volume of which is slightly larger than that of the liquid, most of the container is filled with the liquid, but, immediately above the liquid surface, a vapour phase forms, consisting of molecules that have passed through the liquid surface from liquid to gas; the pressure exerted by that vapour phase is called the vapour (or saturation) pressure. For a pure liquid, this pressure depends only on the temperature, the best-known example being the normal boiling point, which is that temperature at which the vapour pressure is equal to the pressure of the atmosphere. Figure 39 shows vapour pressures for a few common liquids. The vapour pressure is one atmosphere at $100^\circ C$ for water, at $78.5^\circ C$ for ethyl alcohol, and at $125.7^\circ C$ for octane. In a liquid solution, the component with the higher vapour pressure is called the light component, and that with the lower vapour pressure is called the heavy component.

In a liquid mixture, the vapour pressure depends not only on the temperature but also on the composition, and the key problem in understanding the properties of solutions lies in determining this composition dependence. The simplest approximation is to assume that, at constant temperature, the vapour pressure of a solution is a linear function of its composition (*i.e.*, as one increases, so does the other in such proportion that, when the values are plotted, the resulting graph is a straight line). A mixture following this approximation is called an ideal solution.

Fugacity. In a pure liquid, the vapour generated by its escaping molecules necessarily has the same composition as that of the liquid. In a mixture, however, the compo-

Role of atomic weight

Formula weight

Vapour pressure

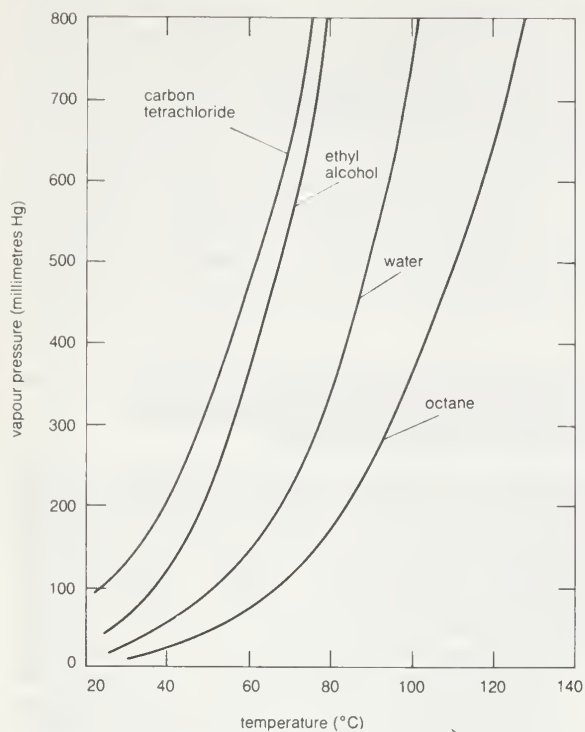


Figure 39: Vapour pressures of several pure liquids.

sition of the vapour is not the same as that of the liquid; the vapour is richer in that component whose molecules have greater tendency to escape from the liquid phase. This tendency is measured by fugacity, a term derived from the Latin *figere* ("to escape, to fly away"). The fugacity of a component in a mixture is (essentially) the pressure that the component exerts in the vapour phase when the vapour is in equilibrium with the liquid mixture. (A state of equilibrium is attained when all the properties remain constant in time and there is no net transfer of energy or matter between the vapour and the liquid.) If the vapour phase can be considered to be an ideal gas (*i.e.*, the molecules in the gas phase are assumed to act independently and without any influence on each other), then the fugacity of a component, i , is equal to its partial pressure, which is defined as the product of the total vapour pressure, P , and the vapour-phase mole fraction, y_i . Assuming ideal gas behaviour for the vapour phase, the fugacity ($y_i P$) equals the product of the liquid-phase mole fraction, x_i , the vapour pressure of pure liquid at the same temperature as that of the mixture, P_i° , and the activity coefficient, γ_i . The real concentration of a substance may not be an accurate measure of its effectiveness, because of physical and chemical interactions, in which case an effective concentration must be used, called the activity. The activity is given by the product of the mole fraction x_i and the activity coefficient γ_i . The equation is:

$$y_i P = \gamma_i x_i P_i^\circ \quad (4)$$

Raoult's law. In a real solution, the activity coefficient, γ_i , depends on both temperature and composition, but, in an ideal solution, γ_i equals 1 for all components in the mixture. For an ideal binary mixture then, the above equation becomes, for components 1 and 2, $y_1 P = x_1 P_1^\circ$ and $y_2 P = x_2 P_2^\circ$, respectively. Upon adding these equations—recalling that $x_1 + x_2 = 1$ and $y_1 + y_2 = 1$ —the total pressure, P , is shown to be expressed by the equation $P = x_1 P_1^\circ + x_2 P_2^\circ = x_1 [P_1^\circ - P_2^\circ] + P_2^\circ$, which is a linear function of x_1 .

Assuming $\gamma_1 = \gamma_2 = 1$, equations for $y_1 P$ and $y_2 P$ express what is commonly known as Raoult's law, which states that at constant temperature the partial pressure of a component in a liquid mixture is proportional to its mole fraction in that mixture (*i.e.*, each component exerts a pressure that depends directly on the number of its molecules present). It is unfortunate that the word law is associated with this relation, because only very few mix-

tures behave according to the equations for ideal binary mixtures. In most cases the activity coefficient, γ_i , is not equal to unity. When γ_i is greater than 1, there are positive deviations from Raoult's law; when γ_i is less than 1, there are negative deviations from Raoult's law.

An example of a binary system that exhibits positive deviations from Raoult's law is represented in Figure 40, the partial pressures and the total pressure being related to the liquid-phase composition: if Raoult's law were valid, all the lines would be straight, as indicated by the dashed lines. As a practical result of these relationships, it is often possible by a series of repeated vaporizations and condensations to separate a liquid mixture into its components, a sequence of steps called fractional distillation.

When the vapour in equilibrium with a liquid mixture has a composition identical to that of the liquid, the mixture is called an azeotrope. It is not possible to separate an azeotropic mixture by fractional distillation because no change in composition is achieved by a series of vaporizations and condensations. Azeotropic mixtures are common. At the azeotropic composition, the total pressure (at constant temperature) is always either a maximum or a minimum with respect to composition, and the boiling temperature (at constant pressure) is always either a minimum or a maximum temperature.

Partial miscibility. Only pairs of liquids that are completely miscible have been considered so far. Many pairs of liquids, however, are only partially miscible in one another, the degree of miscibility often depending strongly on temperature. In most cases, rising temperature produces enhanced solubility, but this is not always so. For example, at 50° C the solubility (weight percent) of *n*-butyl alcohol in water is 6.5 percent, whereas that of water in *n*-butyl alcohol is 22.4 percent. At 127° C, the upper consolute temperature, complete miscibility is attained: above 127° C the two liquids mix in all proportions, but below 127° C they show a miscibility gap. Thus, if *n*-butyl alcohol is added to water at 50° C, there is only one liquid phase until 6.5 weight percent of the mixture is alcohol; when more alcohol is added, a second liquid phase appears the composition of which is 22.4 weight percent water. When sufficient alcohol is present to make the overall compo-

Azeotropic mixtures

The activity coefficient

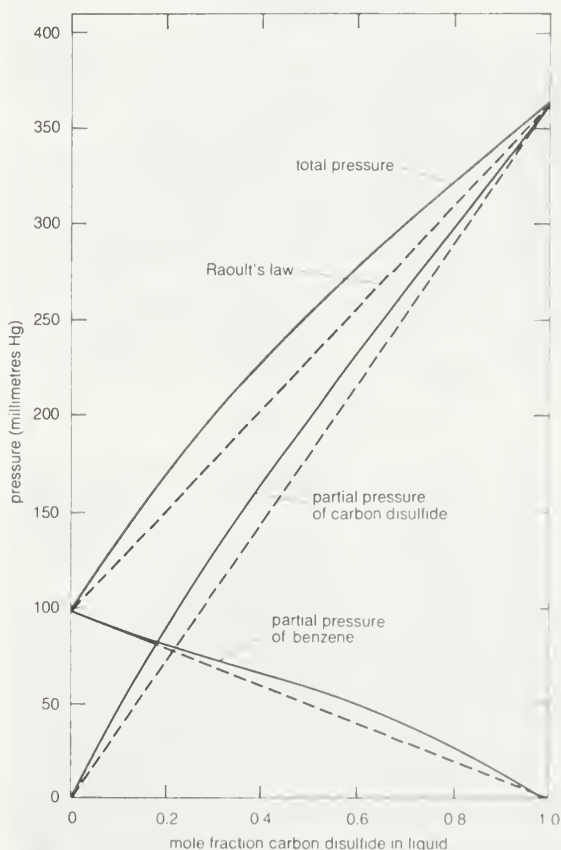


Figure 40: Total pressure and partial pressures for the system benzene-carbon disulfide at 25° C (see text).

sition 77.6 weight percent alcohol, the first phase disappears, and only one liquid phase remains. A qualitatively different example is the system water-triethylamine, which has a lower consolute temperature at 17° C. Below 17° C the two liquids are completely miscible, but at higher temperatures they are only partially miscible. Finally, it is possible, although rare, for a binary system to exhibit both upper and lower consolute temperatures. Above 128° C and below 49° C butyl glycol and water are completely miscible, but between these temperatures they do not mix in all proportions.

Colligative properties. Colligative properties depend only on the concentration of the solute, not on the identity of the solute molecules. The concept of an ideal solution, as expressed by Raoult's law, was already well-known during the last quarter of the 19th century, and it provided the early physical chemists with a powerful technique for measuring molecular weights. (Reliable measurements of molecular weights, in turn, provided important evidence for the modern atomic and molecular theory of matter.)

Rise in boiling point. It was observed that, whenever one component in a binary solution is present in large excess, the partial pressure of that component is correctly predicted by Raoult's law, even though the solution may exhibit departures from ideal behaviour in other respects. When Raoult's law is applied to the solvent of a very dilute solution containing a nonvolatile solute, it is possible to calculate the mole fraction of the solute from an experimental determination of the rise in boiling point that results when the solute is dissolved in the solvent. Since the separate weights of solute and solvent are readily measured, the procedure provides a simple experimental method for the determination of molecular weight. If a weighed amount of a nonvolatile substance, w_2 , is dissolved in a weighed amount of a solvent, w_1 , at constant pressure, the increase in the boiling temperature, ΔT_{b1} , the gas constant, R (derived from the gas laws), the heat of vaporization of the pure solvent per unit weight, l_1^{vap} , and the boiling temperature of pure solvent, T_{b1} , are related in a simple product of ratios equal to the molecular weight of the solute, M_2 . The equation is:

$$M_2 = \left(\frac{RT_{b1}^2}{l_1^{\text{vap}}} \right) \left(\frac{w_2}{w_1} \right) \left(\frac{1}{\Delta T_{b1}} \right). \quad (5)$$

The essence of this technique follows from the observation that, in a dilute solution of a nonvolatile solute, the rise in boiling point is proportional to the number of solute molecules, regardless of their size and mass.

Decrease in freezing point. Another colligative property of solutions is the decrease in the freezing temperature of a solvent that is observed when a small amount of solute is dissolved in that solvent. By reasoning similar to that leading to equation (5), the freezing-point depression, ΔT_f , the freezing temperature of pure solvent, T_{f1} , the heat of fusion (also called the heat of melting) of pure solvent per unit weight, l_1^{fusion} , and the weights of solute and solvent in the solution, w_2 and w_1 , respectively, are so related as to equal the molecular weight of solute, M_2 , in the equation

$$M_2 = \left(\frac{RT_{f1}^2}{l_1^{\text{fusion}}} \right) \left(\frac{w_2}{w_1} \right) \left(\frac{1}{\Delta T_f} \right). \quad (6)$$

A well-known practical application of freezing-point depression is provided by adding antifreeze to the cooling water in an automobile's radiator. Water alone freezes at 0° C, but the freezing temperature decreases appreciably when ethylene glycol is mixed with water.

Osmotic pressure. A third colligative property, osmotic pressure, helped to establish the fundamentals of modern physical chemistry and played a particularly important role in the early days of solution theory. Osmosis is especially important in medicine and biology, but in recent years it has also been applied industrially to problems such as the concentration of fruit juices, the desalting of seawater, and the purification of municipal sewage. Osmosis occurs whenever a liquid solution is in contact with a semipermeable membrane—*i.e.*, a thin, porous wall whose porosity is such that some, but not all, of the components in the liquid mixture can pass through the wall. A semipermeable membrane is a selective barrier, and many

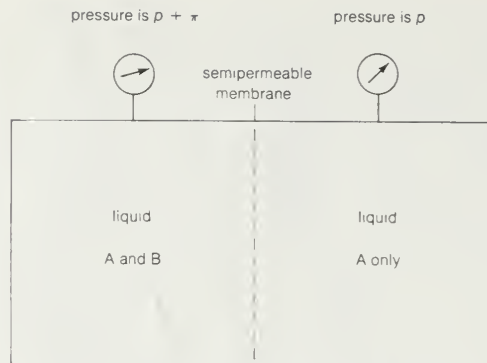


Figure 41: Osmotic pressure π caused by a membrane that allows A to pass but not B. A representative system could consist of water (A) and salt (B).

such barriers are found in plants and animals. Osmosis gives rise to what is known as osmotic pressure, as illustrated in Figure 41, which shows a container at uniform pressure divided into two parts by a semipermeable membrane that allows only molecules of component A to pass from the left to the right side; the selective membrane does not allow molecules of component B to pass. Example compounds for A and B might be water and sodium chloride (table salt), respectively. Molecules of component A are free to pass back and forth through the membrane, but, at equilibrium, when the fugacity (escaping tendency) of A in the right-hand side is the same as that in the left-hand side, there is no net transfer of A from one side to the other. On the left side, the presence of B molecules lowers the fugacity of A, and, therefore, to achieve equal fugacities for A on both sides, some compensating effect is needed on the left side. This compensating effect is an enhanced pressure, designated by π and called osmotic pressure. At equilibrium the pressure in the left side of the container is larger than that in the right side; the difference in pressure is π . In the simplest case, when the concentration of B is small (*i.e.*, A is in excess), the osmotic pressure is the product of the gas constant (R), the absolute temperature (T), and the concentration of B (c_B) in the solution expressed in moles of B per unit volume: $\pi = RTc_B$. Since the osmotic pressure for a dilute solution is proportional to the number of solute molecules, it is a colligative property, and, as a result, osmotic-pressure measurements are often used to determine molecular weights, especially for large molecules such as polymers. When w_B grams of solute B are added to a large amount of solvent A at temperature T , and V is the volume of liquid solvent A in the left side of the container, then the molecular weight of B, M_B , is given by

$$M_B = \frac{w_B RT}{\pi V}. \quad (7)$$

For sodium chloride in water, c_B is the concentration of the ions, which is twice the concentration of the salt owing to the dissociation of the salt (NaCl) into sodium ions (Na^+) and chloride ions (Cl^-). Thus, for a 3.5 percent sodium chloride solution at 25° C, π is 29 atmospheres, which is the minimum pressure at which a desalination reverse osmosis process can operate.

Transport properties in solutions. Pure fluids have two transport properties that are of primary importance: viscosity and thermal conductivity. Transport properties differ from equilibrium properties in that they reflect not what happens at equilibrium but the speed at which equilibrium is attained. In solutions these two transport properties are also important. In addition, there is a third one, called diffusivity.

Viscosity. The viscosity of a fluid (pure or not) is a measure of its ability to resist deformation. If water is poured into a thin vertical tube with a funnel at the top, it flows easily through the tube, but salad oil is difficult to force into the tube. If the oil is heated, however, its flow through the tube is much facilitated. The intrinsic property that is responsible for these phenomena is the viscosity (the "thickness") of the fluid, a property which is often

Osmotic equilibrium

strongly affected by temperature. All fluids (liquid or gas) exhibit viscous behaviour (*i.e.*, all fluids resist deformation to some degree), but the range of viscosity is enormous: the viscosity of air is extremely small, while that of glass is essentially infinite. The viscosity of a solution depends not only on temperature but also on composition. By varying the composition of a petroleum mixture, it is possible to attain a desired viscosity at a particular temperature. This is precisely what the oil companies do when they sell oil to a motorist: in winter, they recommend an oil with lower viscosity than that used in summer, because otherwise, on a cold morning, the viscosity of the lubricating oil may be so high that the car's battery will not be powerful enough to move the lubricated piston.

Thermal conductivity. The thermal conductivity of a material reflects its ability to transfer heat by conduction. In practical situations both viscosity and thermal conductivity are important, as is illustrated by the contrast between an air mattress and a water bed. Because of its low viscosity, air yields rapidly to an imposed load, and thus the air mattress responds quickly when someone lying on it changes position. Water, because of its higher viscosity, noticeably resists deformation, and someone lying on a water bed experiences a caressing response whenever position is changed. At the same time, since the thermal conductivity as well as the viscosity of water are larger than those of air, the user of a water bed rapidly gets cold unless a heater keeps the water warm. No heater is required by the user of an air mattress because stagnant air is inefficient in removing heat from a warm body.

Composition and temperature affect the thermal conductivity of a solution but, in typical liquid mixtures, the effect on viscosity is much larger than that on thermal conductivity.

Diffusivity. While viscosity is concerned with the transfer of momentum and thermal conductivity with the transfer of heat, diffusivity is concerned with the transport of molecules in a mixture. If a lump of sugar is put into a cup of coffee, the sugar molecules travel from the surface of the lump into the coffee at a speed determined by the temperature and by the pertinent intermolecular forces. The characteristic property that determines this speed is called diffusivity—*i.e.*, the ability of a molecule to diffuse through a sea of other molecules. Diffusivities in solids are extremely small, and those in liquids are much smaller than those in gases. For this reason, a spoon is used to stir the coffee to speed up the motion of the sugar molecules, but, if the odour of cigarette smoke fills a room, little effort is needed to clear the air—opening the windows for a few minutes is sufficient.

In order to define diffusivity, it is necessary to consider a binary fluid mixture in which the concentration of solute molecules is c_1 at position 1 and c_2 at position 2, which is l centimetres from position 1; if c_1 is larger than c_2 , then the concentration gradient (change with respect to distance), given by $(c_2 - c_1)/l$, is a negative number, indicating that molecules of solute spontaneously diffuse from position 1 to position 2. The number of solute molecules that pass through an area of one square centimetre perpendicular to l , per second, is called the flux J (expressed in molecules per second per square centimetre). The diffusivity D is given by the formula

$$D = - \frac{J}{(c_2 - c_1)/l} \quad (8)$$

The leading minus sign is introduced because, when the gradient is positive, J is negative, and, by convention, D is a positive number. In binary gaseous mixtures, diffusivity depends only weakly on the composition, and, therefore, to a good approximation, the diffusivity of gas A in gas B is the same as that of gas B in gas A. In liquid systems, however, the diffusivity of solute A in solvent B may be significantly different from that of solute B in solvent A. In a very viscous fluid, molecules cannot rapidly move from one place to another. Therefore, in liquid systems, the diffusivity of solute A depends strongly on the viscosity of solvent B and vice versa. While the letter D is always used for diffusivity, viscosity is commonly given the symbol η ; in many liquid solutions it is observed that, as the

composition changes (as long as the temperature remains constant), the product $D\eta$ remains nearly the same.

OTHER THERMODYNAMICS AND INTERMOLECULAR FORCES IN SOLUTIONS

The properties of solutions depend, essentially, on two characteristics: first, the manner in which the molecules arrange themselves (that is, the geometric array in which the molecules occupy space) and, second, the nature and strength of the forces operating between the molecules.

Energy considerations. The first characteristic is reflected primarily in the thermodynamic quantity S , called entropy, which is a measure of disorder, and the second characteristic is reflected in the thermodynamic quantity H , called enthalpy, which is a measure of potential energy—*i.e.*, the energy that must be supplied to separate all the molecules from one another. Enthalpy minus the product of the absolute temperature T and entropy equals a thermodynamic quantity G , called Gibbs energy (also called free energy):

$$G = H - TS \quad (9)$$

From the second law of thermodynamics, it can be shown that, at constant temperature and pressure, any spontaneous process is accompanied by a decrease in Gibbs energy. The change in G that results from mixing is designated by ΔG , which, in turn, is related to changes in H and S at constant temperature by the equation

$$\Delta G = \Delta H - T\Delta S \quad (10)$$

At a fixed temperature and pressure, two substances mix spontaneously whenever ΔG is negative; that is, mixing (either partial or complete) occurs whenever the Gibbs energy of the substances after mixing is less than that before mixing.

The two characteristics that determine solution behaviour, structure and intermolecular forces, are, unfortunately, not independent, because the structure is influenced by the intermolecular forces and because the potential energy of the mixture depends on the structure. Only in limiting cases is it possible, on the one hand, to calculate ΔS (the entropy change upon mixing) from structural considerations alone and, on the other, to calculate ΔH (the enthalpy change of mixing) exclusively from relations describing intermolecular forces. Nevertheless, such calculations have proved to be useful for establishing models that approximate solution behaviour and that serve as guides in interpreting experimental measurements. Solutions for which structural considerations are dominant are called athermal solutions, and those for which the effects of intermolecular forces are more important than those of structure are called regular solutions (see below *Theories of solutions: Solutions of non-electrolytes: Regular solutions and Athermal solutions*).

Effects of molecular structure. A variety of forces operate between molecules, and there is a qualitative relation between the properties of a solution and the types of intermolecular forces that operate within it. The volume occupied by a solution is determined primarily by repulsive forces. When two molecules are extremely close to one another, they must necessarily exert a repulsive force on each other since two molecules of finite dimensions cannot occupy the same space; two molecules in very close proximity resist attempts to shorten the distance between them.

At larger distances of separation, molecules may attract or repel each other depending on the sign (plus or minus) and distribution of their electrical charge. Two ions attract one another if the charge on one is positive and that on the other is negative; they repel when both carry charges of the same sign. Forces between ions are called Coulomb forces and are characterized by their long range; the force (F) between two ions is inversely proportional to the square of the distance between them; *i.e.*, F varies as $1/r^2$. Noncoulombic physical forces between molecules decay more rapidly with distance; *i.e.*, in general F varies as $1/r^n$, n being larger than 2 for intermolecular forces other than those between ions.

The Coulomb force (F) equals the product of the magnitude of the charge on one ion (e_1) and that on the other

(e_2) divided by the product of the distance squared (r^2) and the dielectric constant (ϵ):

$$F = \frac{e_1 e_2}{r^2 \epsilon} \quad (11)$$

If both e_1 and e_2 are positive, F is positive and the force is repulsive. If either e_1 or e_2 is positive while the other is negative, F is negative and the force is attractive. Coulomb forces are dominant in electrolyte solutions.

Molecular structure and charge distribution. If a molecule has no net electrical charge, its negative charge is equal to its positive charge. The forces experienced by such molecules depend on how the positive and negative charges are arranged in space. If the arrangement is spherically symmetric, the molecule is said to be nonpolar; if there is an excess of positive charge on one end of the molecule and an excess of negative charge on the other, the molecule has a dipole moment (*i.e.*, a measurable tendency to rotate in an electric or magnetic field) and is therefore called polar. The dipole moment (μ) is defined as the product of the magnitude of the charge, e , and the distance separating the positive and negative charges, l : $\mu = el$. Electrical charge is measured in electrostatic units (esu), and the typical charge at one end of a molecule is of the order of 10^{-10} esu; the distance between charges is of the order of 10^{-8} centimetres (cm). Dipole moments, therefore, usually are measured in debyes (one debye is 10^{-18} esu-cm). For nonpolar molecules, $\mu = 0$.

Polar molecules. The force F between two polar molecules is directly proportional to the product of the two dipole moments (μ_1 and μ_2) and inversely proportional to the fourth power of the distance between them (r^4); that is, F varies as $\mu_1 \mu_2 / r^4$. The equation for this relationship contains a constant of proportionality ($F = k \mu_1 \mu_2 / r^4$), the sign and magnitude of which depend on the mutual orientation of the two dipoles; if the positive end of one faces the negative end of the other, the constant of proportionality is negative (meaning that an attractive force exists), while it is positive (meaning that a repulsive force exists) when the positive end of one faces the positive end of the other. When polar molecules are free to rotate, they tend to favour those orientations that lead to attractive forces. To a first approximation, the force (averaged over all orientations) is inversely proportional to the temperature and to the seventh power of the distance of separation. Mixtures of polar molecules often exhibit only mild deviations from ideality, but mixtures containing polar and nonpolar molecules are frequently strongly nonideal. Because of the qualitative and quantitative differences in intermolecular forces, the molecules segregate: the polar molecules prefer to be with each other, and so do the nonpolar ones. Only at higher temperatures, such that the thermal energy of the molecules offsets the cohesion between identical molecules, do the two liquids mix in all proportions. In mixtures containing both polar and nonpolar components, deviations from Raoult's law diminish as temperature rises.

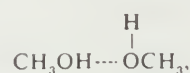
Nonpolar molecules. A nonpolar molecule is one whose charge distribution is spherically symmetric when averaged over time; since the charges oscillate, a temporary dipole moment exists at any given instant in a so-called nonpolar molecule. These temporary dipole moments fluctuate rapidly in magnitude and direction, giving rise to intermolecular forces of attraction called London (or dispersion) forces. All molecules, charged or not, polar or not, interact by London forces. To a first approximation, the London force between two molecules is inversely proportional to the seventh power of the distance of separation; it is therefore short-range, decreasing rapidly as one molecule moves away from the other. The London theory indicates that for simple molecules positive deviations from Raoult's law may be expected (*i.e.*, the activity coefficient γ_i is greater than 1, as explained previously). Since the London theory suggests that the attractive forces between unlike simple molecules are smaller than those corresponding to an ideal solution, the escaping tendency of the molecules in solution is larger than that calculated by Raoult's law. As a result, mixing of small nonpolar molecules is endothermic (absorbing heat from the sur-

roundings) and the volume occupied by the liquid solution often exceeds that of the unmixed components—that is, the components expand on mixing.

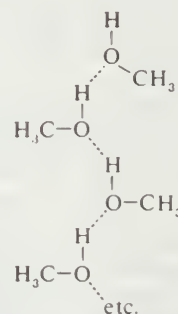
In addition to the forces listed above, there are so-called induction forces set up when a charged or polar molecule induces a dipole in another molecule: the electric field of the inducing molecule distorts the charge distribution in the other. When a charged molecule induces a dipole in another, the force is always attractive and is inversely proportional to the fifth power of the distance of separation. When a polar molecule induces a dipole in another molecule, the force is also attractive and is inversely proportional to the seventh power of separation. Induction forces are usually small but may make a significant contribution to the energy of a mixture of molecules that are strongly dissimilar.

Effects of chemical interactions. In many cases the properties of a mixture are determined primarily by forces that are more properly classified as chemical rather than as physical. For example, when dinitrogen pentoxide is dissolved in water, a new substance, nitric acid, is formed; and it is necessary to interpret the behaviour of such a solution in terms of its chemical properties, which, in this case, are more important than its physical properties. This example is an extreme one, and there are many solutions for which the chemical effect is less severe but nevertheless dominant.

Hydrogen bonding: association. This dominance is especially important in those solutions that involve hydrogen bonding. Whenever a solution contains molecules with an electropositive hydrogen atom and with an electronegative atom (such as nitrogen, oxygen, sulfur, or fluorine), hydrogen bonding may occur and, when it does, the properties of the solution are affected profoundly. Hydrogen bonds may form between identical molecules or between dissimilar molecules; for example, methanol (CH_3OH) has an electropositive (electron-attracting) hydrogen atom and also an electronegative (electron-donating) oxygen atom, and therefore two methanol molecules may hydrogen-bond (represented by the dashed line) singly to form the structure



or in chains to form



Hydrogen bonding between identical molecules is often called association.

Hydrogen bonding: solvation. In a mixture of methanol and, say, pyridine ($\text{C}_5\text{H}_5\text{N}$), hydrogen bonds can also form between the electropositive hydrogen atom in methanol and the electronegative nitrogen atom in pyridine. Hydrogen bonding between dissimilar molecules is an example of a type of interaction known as solvation. Since the extent of association or solvation or both depends on the concentrations of the solution's components, the partial pressure of a component is not even approximately proportional to its mole fraction as given by Raoult's law; therefore, large deviations from Raoult's law are commonly observed in solutions in which hydrogen bonding is extensive. Broadly speaking, association of one component, but not the other, tends to produce positive deviations from Raoult's law, because the associating component hydrogen-bonds to a smaller extent when it is surrounded by other molecules than it does in the pure state. On the other hand, solvation

Deviations from Raoult's law

Dipole moment

London forces

between dissimilar molecules tends to produce negative deviations from Raoult's law.

THEORIES OF SOLUTIONS

Solutions of nonelectrolytes. *Activity coefficients and excess functions.* As has been explained previously, when actual concentrations do not give simple linear relations for the behaviour of a solution, activity coefficients, symbolized by γ_i , are used in expressing deviations from Raoult's law. Activity coefficients are directly related to excess functions, and, in attempting to understand solution behaviour, it is convenient to characterize nonelectrolyte solutions in terms of these functions. In particular, it is useful to distinguish between two types of limiting behaviour: one corresponds to that of a regular solution; the other, to that of an athermal solution (*i.e.*, when components are mixed, no heat is generated or absorbed).

In a binary mixture with mole fractions x_1 and x_2 and activity coefficients γ_1 and γ_2 , these quantities can be related to a thermodynamic function designated by G^E , called the excess Gibbs (or free) energy. The significance of the word excess lies in the fact that G^E is the Gibbs energy of a solution in excess of what it would be if it were ideal.

In a binary solution the two activity coefficients are not independent but are related by an exact differential equation called the Gibbs-Duhem relation. If experimental data at constant temperature are available for γ_1 and γ_2 as a function of composition, it is possible to apply this equation to check the data for thermodynamic consistency; the data are said to be consistent only if they satisfy the Gibbs-Duhem relation. Experimental data that do not satisfy this relation are thermodynamically inconsistent and therefore must be erroneous.

To establish a theory of solutions, it is necessary to construct a theoretical (or semitheoretical) equation for the excess Gibbs energy as a function of absolute temperature (T) and the mole fractions x_1 and x_2 . After such an equation has been established, the individual activity coefficients can readily be calculated.

Gibbs energy, by definition, consists of two parts: one part is the enthalpy, which reflects the intermolecular forces between the molecules, which, in turn, are responsible for the heat effects that accompany the mixing process (enthalpy is, in a general sense, a measure of the heat content of a substance); and the other part is the entropy, which reflects the state of disorder (a measure of the random behaviour of particles) in the mixture. The excess Gibbs energy G^E is given by the equation

$$G^E = H^E - TS^E, \quad (12)$$

where H^E is the excess enthalpy and S^E is the excess entropy. The word excess means in excess of that which would prevail if the solution were ideal. In the simplest case, both H^E and S^E are zero; in that case the solution is ideal and $\gamma_1 = \gamma_2 = 1$. In the general case, neither H^E nor S^E is zero, but two types of semi-ideal solutions can be designated: in the first, S^E is zero but H^E is not; this is called a regular solution. In the second, H^E is zero but S^E is not; this is called an athermal solution. An ideal solution is both regular and athermal.

Regular solutions. The word regular implies that the molecules mix in a completely random manner, which means that there is no segregation or preference; a given molecule chooses its neighbours with no regard for chemical identity (species 1 or 2). In a regular solution of composition x_1 and x_2 , the probability that the neighbour of a given molecule is of species 1 is given by the mole fraction x_1 , and the probability that it is of species 2 is given by x_2 .

Two liquids form a solution that is approximately regular when the molecules of the two liquids do not differ appreciably in size and there are no strong orienting forces caused by dipoles or hydrogen bonding. In that event, the mixing process can be represented by the lattice model shown in Figure 42: the left half of the diagram shows pure liquids 1 and 2, and the right half shows the mixture obtained when the central molecule of liquid 1 is interchanged with the central molecule of liquid 2. Before interchange, the potential energy between central

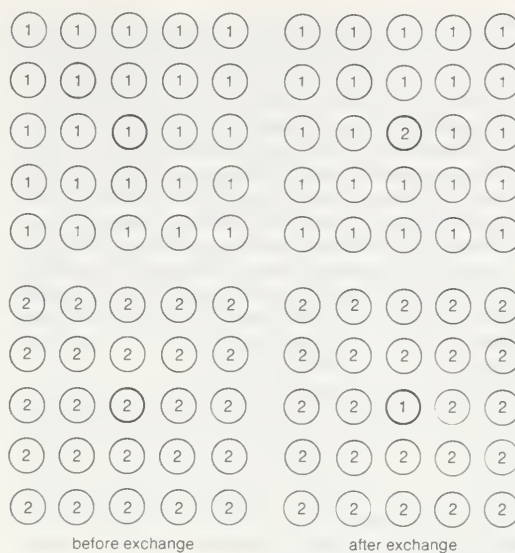


Figure 42: Physical significance of interchange energy. The energy absorbed in the process above is 2ω (see text).

molecule 1 and one of its immediate neighbours is Γ_{11} , and that between central molecule 2 and one of its immediate neighbours is Γ_{22} . After interchange, the potential energy between molecule 1 and one of its immediate neighbours is Γ_{12} , and that between molecule 2 and one of its immediate neighbours is also Γ_{12} . The change in energy that accompanies this mixing process is equal to twice the interchange energy (ω), which is equal to the potential energy after mixing less one-half the sum of the potential energies before mixing, the whole multiplied by the number of immediate neighbours, called the coordination number (z), surrounding the two shifted molecules:

$$\omega = z[\Gamma_{12} - \frac{1}{2}(\Gamma_{11} + \Gamma_{22})]. \quad (13)$$

In the two-dimensional representation (Figure 42), z equals 4; but, in three dimensions, z varies between 6 and 12, depending on the lattice geometry. In this simple lattice model, the interchange process occurs without change of volume; thus, in this particular case, the excess enthalpy is the same as the energy change upon mixing. Assuming regular-solution behaviour (*i.e.*, $S^E = 0$), an equation may be derived relating Gibbs energy, Avogadro's number, interchange energy, and mole fractions. In principle, the interchange energy (ω) may be positive or negative, but, for simple molecules, for which only London forces of attraction are important, ω is positive. The equation obtained from the simple lattice model can be extended semiempirically to apply to mixtures of molecules whose sizes are not nearly the same by using volume fractions instead of mole fractions to express the effect of composition and by introducing the concept of cohesive energy density, which is defined as the potential energy of a liquid divided by its volume. The adjective cohesive is well chosen because it indicates that this energy is associated with the forces that keep the molecules close together in a condensed state. Again restricting attention to nonpolar molecules and assuming a completely random mixture ($S^E = 0$), an equation may be derived that requires only pure-component properties to predict the excess Gibbs energy (and hence the activity coefficients) of binary mixtures. Because of many simplifying assumptions, this equation does not give consistently accurate results, but in many cases it provides good semiquantitative estimates. The form of the equation is such that the excess Gibbs energy is larger than zero; hence, the equation is not applicable to mixtures that have negative deviations from Raoult's law.

Athermal solutions. In a solution in which the molecules of one component are much larger than those of the other, the assumption that the solution is regular (*i.e.*, that $S^E = 0$) no longer provides a reasonable approximation even if the effect of intermolecular forces is neglected. A large flexible molecule (*e.g.*, a chain molecule such as polyethylene) can attain many more configurations when

it is surrounded by small molecules than it can when surrounded by other large flexible molecules; the state of disorder in such a solution is therefore much larger than that of a regular solution in which $S^E = 0$. A solution of very large molecules (*i.e.*, polymers) in an ordinary liquid solvent is analogous to a mixture of cooked spaghetti (representing the polymers) and tomato sauce (the solvent). When there is a large amount of sauce and relatively little spaghetti, each piece of spaghetti is free to exist in many different shapes; this freedom, however, becomes restricted as the number of spaghetti pieces rises and the amount of sauce available for each strand declines. The excess entropy then is determined primarily by the freedom that the spaghetti has in the tomato sauce mixture relative to the freedom it has in the absence of sauce.

Regular solutions and athermal solutions represent limiting cases; real solutions are neither regular nor athermal. For real solutions it has been proposed to calculate G^E by combining the equations derived separately for regular solutions and for athermal solutions, but, in view of the restrictive and mutually inconsistent assumptions that were made in deriving these two equations, the proposal has met with only limited success.

Associated and solvated solutions. For those solutions in which there are strong intermolecular forces due to large dipole moments, hydrogen bonding, or complex formation, equations based on fundamental molecular theory cannot be applied, but it is frequently useful to apply a chemical treatment—*i.e.*, to describe the liquid mixture in terms of association and solvation, by assuming the existence of a variety of distinct chemical species in chemical equilibrium with one another. For example, there is much experimental evidence for association in acetic acid, in which most of the molecules dimerize; *i.e.*, two single acetic acid molecules, called monomers, combine to form a new molecule, called a dimer, through hydrogen bonding. When acetic acid is dissolved in a solvent such as benzene, the extent of dimerization of acetic acid depends on the temperature and on the total concentration of acetic acid in the solution. The escaping tendency (vapour pressure) of a monomer is much greater than that of a dimer, and it is thus possible to explain the variation of activity coefficient with composition for acetic acid in benzene; the activity coefficient of acetic acid in an excess of benzene is large because, under these conditions, acetic acid is primarily in the monomeric state, whereas pure acetic acid is almost completely dimerized. In the acetic acid-benzene system, association of acetic acid molecules produces positive deviations from Raoult's law.

Solvation

When a solvent and a solute molecule link together with weak bonds, the process is called solvation. For example, in the system acetone-chloroform, a hydrogen bond is formed between the hydrogen atom in chloroform and the oxygen atom in acetone. In this case, hydrogen bonding depresses the escaping tendencies of both components, producing negative deviations from Raoult's law.

While hydrogen bonding is frequently encountered in solutions, there are many other examples of weak chemical-bond formation between dissimilar molecules. The formation of such weak bonds is called complex formation—that is, formation of a new chemical species, called a complex, which is held together by weak forces that are chemical in nature rather than physical. Such complexes usually exist only in solution; because of their low stability, they cannot, in general, be isolated. The ability of molecules to form complexes has a strong effect on solution behaviour. For example, the solubility of a sparingly soluble species can be much increased by complex formation: the solubility of silver chloride in water is extremely small since silver chloride dissociates only slightly to silver ion and chloride ion; however, when a small quantity of ammonia is added, solubility rises dramatically because of the reaction of six molecules of ammonia with one silver ion to form the complex ion $\text{Ag}(\text{NH}_3)_6^+$. By tying up silver ions and forcing extensive dissociation of molecular silver chloride, the ammonia pulls silver chloride into aqueous solution.

In recent years there has been much interest in the use of computers to generate theoretical expressions for the activity coefficients of solutions. In many cases the

calculations have been restricted to model systems, in particular to mixtures of hard-sphere (envisioned as billiard balls) molecules—*i.e.*, idealized molecules that have finite size but no forces of attraction. These calculations have produced a better understanding of the structure of simple liquid solutions since the manner in which nonpolar and non-hydrogen-bonding molecules arrange themselves in space is determined primarily by their size and shape and only secondarily by their attractive intermolecular forces. The results obtained for hard-sphere molecules can be extended to real molecules by applying corrections required for attractive forces and for the "softness" of the molecules—*i.e.*, the ability of molecules to interpenetrate (overlap) at high temperatures. While practical results are still severely limited and while the amount of required computer calculation is large even for simple binary systems, there is good reason to believe that advances in the theory of solution will increasingly depend on computerized, as opposed to analytical, models.

Solutions of electrolytes. Near the end of the 19th century, the properties of electrolyte solutions were investigated extensively by the early workers in physical chemistry. A suggestion of Svante August Arrhenius, a Swedish chemist, that salts of strong acids and bases (for example, sodium chloride) are completely dissociated into ions when in aqueous solution received strong support from electrical-conductivity measurements and from molecular-weight studies (freezing-point depression, boiling-point elevation, and osmotic pressure). These studies showed that the number of solute particles was larger than it would be if no dissociation occurred. For example, a 0.001 molal solution of a uni-univalent electrolyte (one in which each ion has a valence, or charge, of 1, and, when dissociated, two ions are produced) such as sodium chloride, Na^+Cl^- , exhibits colligative properties corresponding to a nonelectrolyte solution whose molality is 0.002; the colligative properties of a 0.001 molal solution of a univalent-divalent electrolyte (yielding three ions) such as magnesium bromide, $\text{Mg}^{2+}\text{Br}_2^-$, correspond to those of a nonelectrolyte solution with a molality of 0.003. At somewhat higher concentrations the experimental data showed some inconsistencies with Arrhenius' dissociation theory, and initially these were ascribed to incomplete, or partial, dissociation. In the years 1920–30, however, it was shown that these inconsistencies could be explained by electrostatic interactions (Coulomb forces) of the ions in solution. The current view of electrolyte solutions is that, in water at normal temperatures, the salts of strong acids and strong bases are completely dissociated into ions at all concentrations up to the solubility limit. At high concentrations Coulombic interactions may cause the formation of ion pairs, which implies that the ions are not dispersed uniformly in the solution but have a tendency to form two-ion aggregates in which a positive ion seeks the close proximity of a negative ion and vice versa. While the theory of dilute electrolyte solutions is well advanced, no adequate theory exists for concentrated electrolyte solutions primarily because of the long-range Coulomb forces that dominate in ionic solutions.

The equilibrium properties of electrolyte solutions can be studied experimentally by electrochemical measurements, freezing-point depressions, solubility determinations, osmotic pressures, or measurements of vapour pressure. Most electrolytes, such as salts, are nonvolatile at ordinary temperature, and, in that event, the vapour pressure exerted by the solution is the same as the partial pressure of the solvent. The activity coefficient of the solvent can, therefore, be found from total-pressure measurements, and, using the Gibbs-Duhem equation, it is then possible to calculate the activity coefficient of the electrolyte solute. This activity coefficient is designated by γ_{\pm} to indicate that it is a mean activity coefficient for the positive and negative ions. Since it is impossible to isolate positive ions and negative ions into separate containers, it is not possible to determine individual activity coefficients for the positive ions and for the negative ions. The mean activity coefficient γ_{\pm} is so defined that it approaches a value of unity at very low molality where the ions are so far apart that they exert negligible influence on one another. For

Arrhenius' theory of salts

Current theory of salts

small concentrations of electrolyte, the theory of Peter Debye, a Dutch-born American physical chemist, and Erich Hückel, a German chemist, relates γ_{\pm} to the ionic strength, which is the sum of the products of the concentration of each ion (in moles per litre) and the square of its charge; the equation predicts that γ_{\pm} decreases with rising ionic strength in agreement with experiment at very low ionic strength; at higher ionic strength, however, γ_{\pm} rises, and in some cases γ_{\pm} is greater than 1. The derivation of the Debye-Hückel theory clearly shows that it is limited to low concentrations. Many attempts have been made to extend the Debye-Hückel equation to higher electrolyte concentrations. One of the more successful attempts is based on the idea that the ions are solvated, which means that every ion is surrounded by a tight-fitting shell of solvent molecules.

The concept of solvation is often used to explain properties of aqueous solutions; one well-known property is the salting-out effect, in which the solubility of a nonelectrolyte in water is decreased when electrolyte is added. For example, the solubility of ethyl ether in water at 25° C is 0.91 mole percent, but, in an aqueous solution containing 15 weight percent sodium chloride, it is only 0.13 mole percent. This decrease in solubility can be explained by postulating that some of the water molecules cannot participate in the dissolution of the ether because they are tightly held (solvated) by sodium and chloride ions.

Electrolyte solutions have long been of interest in industry since many common inorganic chemicals are directly obtained, or else separated, by crystallization from aqueous solution. Further, many important chemical and metallurgical products (e.g., aluminum) are obtained or refined by electrochemical processes that occur in liquid solution. In recent years there has been renewed interest in electrolyte solutions because of their relevance to fuel cells as a possible source of power for automobiles.

The properties of electrolyte solutions also have large importance in physiology. Many molecules that occur in biological systems bear electric charges; a large molecule that has a positive electric charge at one end and a negative charge at the other is called a zwitterion. Very large molecules, such as those of proteins, may have numerous positive and negative charges; such molecules are called polyelectrolytes. In solution, the conformation (*i.e.*, the three-dimensional structure) of a large, charged molecule is strongly dependent on the ionic strength of the dissolving medium; for example, depending on the nature and concentration of salts present in the solvent, a polyelectrolyte molecule may coagulate into a ball, it may stretch out like a rod, or it may form a coil or helix. The conformation, in turn, is closely related to the molecule's physiological function. As a result, improved understanding of the properties of electrolyte solutions has direct consequences in molecular biology and medicine.

SOLUBILITIES OF SOLIDS AND GASES

Since the dissolution of one substance in another can occur only if there is a decrease in the Gibbs energy, it follows that, generally speaking, gases and solids do not dissolve in liquids as readily as do other liquids. To understand this, the dissolution of a solid can be visualized as occurring in two steps: in the first, the pure solid is melted at constant temperature to a pure liquid, and, in the second, that liquid is dissolved at constant temperature in the solvent. Similarly, the dissolution of a gas can be divided at some fixed pressure into two parts, the first corresponding to constant-temperature condensation of the pure gas to a liquid and the second to constant-temperature mixing of that liquid with solvent. In many cases, the pure liquids (obtained by melting or by condensation) may be hypothetical (*i.e.*, unstable and, therefore, physically unobtainable), but usually their properties can be estimated by reasonable extrapolations. It is found that the change in Gibbs energy corresponding to the first step is positive and, hence, in opposition to the change needed for dissolution. For example, at -10° C, ice is more stable than water, and, at 110° C and one atmosphere, steam is more stable than water. Therefore, the Gibbs energy of melting ice at -10° C is positive, and the Gibbs energy

of condensing steam at one atmosphere and 110° C is also positive. For the second step, however, the change in Gibbs energy is negative; its magnitude depends on the equilibrium composition of the mixture. Owing to the positive Gibbs energy change that accompanies the first step, there is a barrier that makes it more difficult to dissolve solids and gases as compared with liquids.

For gases at normal pressures, the positive Gibbs energy of condensation increases with rising temperature, but, for solids, the positive Gibbs energy of melting decreases with rising temperature. For example, the change in energy, ΔG , of condensing steam at one atmosphere is larger at 120° C than it is at 110° C, while the change in energy of melting ice at -5° C is smaller than it is at -10° C. Thus, as temperature rises, the barrier becomes larger for gases but lower for solids, and therefore, with few exceptions, the solubility of a solid rises while the solubility of a gas falls as the temperature is raised.

For solids, the positive Gibbs energy "barrier" depends on the melting temperature. If the melting temperature is much higher than the temperature of the solution, the barrier is large, shrinking to zero when the melting temperature and solution temperature become identical.

Table 6: Solubilities of Some Gases* (mole percent)

	heptane	benzene	water
Hydrogen	0.069	0.026	0.0015
Nitrogen	0.12	0.45	0.0012
Methane	0.47	0.21	0.0024
Carbon dioxide	0.77	0.97	0.0608

*At one atmosphere partial pressure, 25° C.

Table 6 gives the solubilities of some common gases, and Table 7 the solubility of (solid) naphthalene in a few typical solvents. These solubilities illustrate the qualitative rule that "like dissolves like"; thus naphthalene, an aromatic hydrocarbon, dissolves more readily in another aromatic hydrocarbon such as benzene than it does in a chlorinated solvent such as carbon tetrachloride or in a hydrogen-bonded solvent such as methyl alcohol. By similar reasoning, the gas methane, a paraffinic hydrocarbon, dissolves more readily in another paraffin such as hexane than it does in water. In all three solvents, the gas hydrogen (which boils at -252.5° C) is less soluble than nitrogen (which boils at a higher temperature, -195.8° C).

Table 7: Solubility of Naphthalene in Various Solvents*

solvent	mole percent naphthalene
Benzene	24.1
Carbon tetrachloride	20.5
Hexane	9.0
Methyl alcohol	1.8
Water	0.0004

*At 20° C.

While exceptions may occur at very high pressures, the solubility of a gas in a liquid generally rises as the pressure of that gas increases. When the pressure of the gas is much larger than the vapour pressure of the solvent, the solubility is often proportional to the pressure. This proportionality is consistent with Henry's law, which states that, if the gas phase is ideal, the solubility x_2 of gas 2 in solvent 1 is equal to the partial pressure (the vapour-phase mole fraction y_2 times the total pressure P —*i.e.*, y_2P) divided by a temperature-dependent constant, $H_{2,1}$ (called Henry's constant), which is determined to a large extent by the intermolecular forces between solute 2 and solvent 1:

$$x_2 = \frac{y_2 P}{H_{2,1}} \quad (14)$$

When the vapour pressure of solvent 1 is small compared with the total pressure, the vapour-phase mole fraction of gas 2 is approximately one, and the solubility of the gas is proportional to the total pressure. (J.M.P./B.E.P.)

Electrolytes in physiology

Dissolution of a solid

Henry's law

GASEOUS STATE

Structure

The remarkable feature of gases is that they appear to have no structure at all. They have neither a definite size nor shape, whereas ordinary solids have both a definite size and a definite shape, and liquids have a definite size, or volume, even though they adapt their shape to that of the container in which they are placed. Gases will completely fill any closed container; their properties depend on the volume of a container but not on its shape.

KINETIC-MOLECULAR PICTURE

Gases nevertheless do have a structure of sorts on a molecular scale. They consist of a vast number of molecules moving chaotically in all directions and colliding with one another and with the walls of their container. Beyond this, there is no structure—the molecules are distributed essentially randomly in space, traveling in arbitrary directions at speeds that are distributed randomly about an average determined by the gas temperature. The pressure exerted by a gas is the result of the innumerable impacts of the molecules on the container walls and appears steady to human senses because so many collisions occur each second on all sections of the walls. More subtle properties such as heat conductivity, viscosity (resistance to flow), and diffusion are attributed to the molecules themselves carrying the mechanical quantities of energy, momentum, and mass, respectively. These are called transport properties, and the rate of transport is dominated by the collisions between molecules, which force their trajectories into tortuous shapes. The molecular collisions are in turn controlled by the forces between the molecules and are described by the laws of mechanics.

Thus, gases are treated as a large collection of tiny particles subject to the laws of physics. Their properties are attributed primarily to the motion of the molecules and can be explained by the kinetic theory of gases. It is not obvious that this should be the case, and for many years a static picture of gases was instead espoused, in which the pressure, for instance, was attributed to repulsive forces between essentially stationary particles pushing on the container walls. How the kinetic-molecular picture finally came to be universally accepted is a fascinating piece of scientific history and is discussed briefly below in the section *Kinetic theory of gases*. Any theory of gas behaviour based on this kinetic model must also be a statistical one because of the enormous numbers of particles involved. The kinetic theory of gases is now a classical part of statistical physics and is indeed a sort of miniature display case for many of the fundamental concepts and methods of science. Such important modern concepts as distribution functions, cross sections, microscopic reversibility, and time-reversal invariance have their historical roots in kinetic theory, as does the entire atomistic view of matter.

NUMERICAL MAGNITUDES

When considering various physical phenomena, it is helpful for one to have some idea of the numerical magnitudes involved. In particular, there are several characteristics whose values should be known, at least within an order of magnitude (a factor of 10), in order for one to obtain a clear idea of the nature of gaseous molecules. These features include the size, average speed, and intermolecular separation at ordinary temperatures and pressures. In addition, other important considerations are how many collisions a typical molecule makes in one second under these conditions and how far such a typical molecule travels before colliding with another molecule. It has been established that molecules have sizes on the order of a few angstrom units ($1 \text{ \AA} = 10^{-8}$ centimetre [cm]) and that there are about 6×10^{23} molecules in one mole, which is defined as the amount of a substance whose mass in grams is equal to its molecular weight (e.g., 1 mole of water, H_2O , is 18.0152 grams). With this knowledge, one could calculate at least some of the gas values. It is in-

teresting to see how the answers could be estimated from simple observations and then to compare the results to the accepted values that are based on more precise measurements and theories.

One of the easiest properties to work out is the average distance between molecules compared to their diameter; water will be used here for this purpose. Consider 1 gram of H_2O at 100°C and atmospheric pressure, which are the normal boiling point conditions. The liquid occupies a volume of 1.04 cubic centimetres (cm^3); once converted to steam it occupies a volume of $1.67 \times 10^3 \text{ cm}^3$. Thus, the average volume occupied by one molecule in the gas is larger than the corresponding volume occupied in the liquid by a factor of $1.67 \times 10^3 / 1.04$, or about 1,600. Since volume varies as the cube of distance, the ratio of the mean separation distance in the gas to that in the liquid is roughly equal to the cube root of 1,600, or about 12. If the molecules in the liquid are considered to be touching each other, the ratio of the intermolecular separation to the molecular diameter in ordinary gases is on the order of 10 under ordinary conditions. It should be noted that the actual separation and diameter cannot be determined in this way; only their ratio can be calculated.

It is also relatively simple to estimate the average speed of gas molecules. Consider a sound wave in a gas, which is just the propagation of a small pressure disturbance. If pressure is attributed to molecular impacts on a test surface, then surely a pressure disturbance cannot travel faster than the molecules themselves. In other words, the average molecular speed in a gas should be somewhat greater than the speed of sound in the gas. The speed of sound in air at ordinary temperatures is about 330 metres per second (m/s), so the molecular speed will be estimated here to be somewhat greater, say, about 5×10^4 centimetres per second (cm/s). This value depends on the particular gas and the temperature, but it will be sufficient for the kind of estimates sought here.

The average molecular speed, along with an observed rate of the diffusion of gases, can be used to estimate the length and tortuosity of the path traveled by a typical molecule. If a bottle of ammonia is opened in a closed room, at least a few minutes pass before the ammonia can be detected at a distance of just one metre. (Ammonia, NH_3 , is a gas; the familiar bottle of "ammonia" typically seen is actually a solution of the gas in water.) Yet, if the ammonia molecules traveled directly to an observer at a speed somewhat faster than that of sound, the odour should be detectable in only a few milliseconds. The explanation for the discrepancy is that the ammonia molecules collide with many air molecules, and their paths are greatly distorted as a result. For a quantitative estimate of the diffusion time, a more controlled system must be considered, because even gentle stray air currents in a closed room greatly speed up the spreading of the ammonia. To eliminate the effect of such air currents, a closed tube—say, a glass tube one centimetre in diameter and one metre in length—can be used. A small amount of ammonia gas is released at one end, and both ends are then closed. In order to measure how long it takes for the ammonia to travel to the other end, a piece of moist red litmus paper might be used as a detector; it will turn blue when the ammonia reaches it. This process takes quite a long time—about several hours—because diffusion occurs at such a slow rate. In this case, the time will be taken to be approximately 3 hours, or roughly 10^4 seconds (s). During this time interval, a typical ammonia molecule actually travels a distance of $(5 \times 10^4 \text{ cm/s})(10^4 \text{ s}) = 5 \times 10^8 \text{ cm} = 5,000$ kilometres (km), roughly the distance across the United States. In other words, such a molecule travels a total distance of five million metres in order to progress a net distance of only one metre.

The solution to a basic statistical problem can be used to estimate the number of collisions such a typical diffusing molecule experienced (N) and the average distance traveled between collisions (l), called the mean free path.

Random
nature of
gases

Average
molecular
speed

The product of N and l must equal the total distance travelled—*i.e.*, $Nl = 5 \times 10^8$ cm. This distance can be thought of as a chain 5,000 km long, made up of N links, each of length l . The statistical question then is as follows: If such a chain is randomly jumbled, how far apart will its ends be on the average? This end-to-end distance corresponds to the length of the diffusion tube (one metre). This is a venerable statistical problem that recurs in many applications. One of the more vivid ways of illustrating the concept is known as the “drunkard’s walk.” In this scenario a drunkard takes steps of length l but, because of inebriation, takes them in random directions. After N steps, how far will he be from his starting point? The answer is that his progress is proportional not to N but to $N^{1/2}$. For example, if the drunkard takes four steps, each of length l , he will end up at a distance of $2l$ from his starting point. Gas molecules move in three dimensions, whereas the drunkard moves in two dimensions; however, the result is the same. Thus, the square root of N multiplied by the length of the mean free path equals the length of the diffusion tube: $N^{1/2}l = 10^2$ cm. From the equations for Nl and $N^{1/2}l$, it can readily be calculated that $N = 2.5 \times 10^{13}$ collisions and $l = 2.0 \times 10^{-5}$ cm. The mean time between collisions, τ , is found by dividing the time of the diffusion experiment by the number of collisions during that time: $\tau = (10^4)/(2.5 \times 10^{13}) = 4 \times 10^{-10}$ seconds between collisions, corresponding to a collision frequency of 2.5×10^9 collisions per second. It is thus understandable that gases appear to be continuous fluids on ordinary scales of time and distance.

Molecular sizes can be estimated from the foregoing information on the intermolecular separation, speed, mean free path, and collision rate of gas molecules. It would seem logical that large molecules should have a better chance of colliding than do small molecules. The collision frequency and mean free path must therefore be related to molecular size. To find this relationship, consider a single molecule in motion; during a time interval t it will sweep out a certain volume, hitting any other molecules present in this so-called collision volume. If molecules are located by their centres and each molecule has a diameter d , then the collision volume will be a long cylinder of cross-sectional area πd^2 . The cylinder must be sufficiently long to include enough molecules so that good statistics on the number of collisions are obtained, but otherwise the length does not matter. If the molecule is observed for a time t , then the length of the collision cylinder will be $\bar{v}t$, where \bar{v} is the average speed of the molecule, and the volume of the cylinder will be $(\pi d^2)(\bar{v}t)$, the product of its cross-sectional area and its length. Every molecule in the cylinder will be struck within time t , so the number of molecules in the collision cylinder will equal the number of collisions that occur in time t . Each collision will put a kink in the cylinder, but this will not affect the results as long as the number of collisions is not too large. If the gas is uniform, the number of molecules per volume will be consistent throughout the entire gas. Suppose that there are N molecules in volume V ; then there will be $(N/V)(\pi d^2)(\bar{v}t)$ molecules in the collision volume; this is the number of collisions in time t . The mean free path is equal to the total length of the collision cylinder divided by the number of collisions that occur in it:

$$l = \frac{(\bar{v}t)}{(N/V)(\pi d^2)(\bar{v}t)} = \frac{1}{(N/V)(\pi d^2)}$$

Since l has been shown to be roughly 2.0×10^{-5} cm, d could be calculated if N/V was known.

It is relatively easy to find $(N/V)d^3$, from which both d and N/V can be determined. Recall that the volume of one gram of steam is about 1,600 times larger than the volume of one gram of liquid water. In other words, there are roughly 1,600 N molecules in a volume V of liquid, and, if the molecules are just touching (*i.e.*, the separation distance between their centres is one molecular diameter), the volume V of the liquid is $1,600 Nd^3$. When this equation for volume is combined with the above expression for l , the following values are obtained: $d = \pi(2.0 \times 10^{-5})/1,600 = 3.9 \times 10^{-8}$ cm = 3.9 Å, and $N/V = 1/\pi d^2 l = 1.0 \times 10^{19}$ molecules per cubic centimetre.

Thus, a typical molecule is exceedingly small, and there is an impressively large number of them in one cubic centimetre of gas.

Between collisions, a gas molecule travels a distance of about $l/d = (2.0 \times 10^{-5})/(3.9 \times 10^{-8}) = 500$ times its diameter. Since it was calculated above that the average separation between molecules is about 10 times the molecular diameter, the mean free path is approximately 50 times greater than the mean molecular separation. Accordingly, a typical molecule passes roughly 50 other molecules before it hits one.

The following is a summary of the above estimates of molecular quantities in a gas, with a little spread in the numbers to allow for molecules both smaller and larger than the typical ones used here—which are H_2O , NH_3 , and the nitrogen (N_2) plus oxygen (O_2) mixture that is air—and to allow for the fact that some of these quantities depend on temperature and pressure. It is important to note that these estimates and calculations are rather simplified, although fundamentally correct, and that there may well be missing factors such as $3\pi/8$ or $\sqrt{2}$. The numerical estimates for gases at ordinary pressure and temperature are:

molecular diameter	10^{-8} to 10^{-7} cm
molecular number density	10^{19} molecules/cm ³
average molecular speed	10^4 to 10^5 cm/s
average distance between molecules	10^{-7} to 10^{-6} cm
collision rate per molecule	10^9 to 10^{10} collisions/s
average time between collisions	10^{-10} to 10^{-9} s
average distance traveled between collisions (mean free path)	10^{-5} to 10^{-4} cm

The general impression of gas molecules given by these numbers is that they are exceedingly small, that there are enormous numbers of them in even one cubic centimetre, that they are moving very fast, and that they collide many times in one second. Two other facts are especially important. The first is that the lengths involved, especially the mean free path, are minute compared with ordinary lengths, even with the diameter of a capillary tube. This means that gas behaviour and properties are dominated by collisions between molecules and that collisions with walls play only a secondary (though important) role. The second is that the mean free path is much larger than the molecular diameter. Thus, collisions between pairs of molecules are of paramount importance in determining ordinary gas behaviour, while collisions that involve three or more molecules at the same time can basically be ignored.

A cautious reader might feel a bit uneasy about the glibness of the preceding estimates, so a simple check will be made here by calculating the number of molecules in one mole of gas, a quantity known as Avogadro’s number. The number density of a gas was approximated to be about 1.0×10^{19} molecules per cubic centimetre, and from experiment it is known that 1 mole of gas occupies a volume of about 25 litres (2.5×10^4 cubic centimetres) under ordinary conditions. Using these values, an estimate of Avogadro’s number is $(1.0 \times 10^{19})(2.5 \times 10^4) = 2.5 \times 10^{23}$ molecules per mole. This deviates somewhat from the accepted value of 6.022×10^{23} molecules per mole, but the order of magnitude is certainly correct. In point of historical fact, a value for Avogadro’s number as good as this estimate was not obtained until 1865, when Josef Loschmidt in Vienna made a calculation similar to the one here but based on gas viscosity rather than on gas diffusion. In the older German scientific literature, Avogadro’s number is often referred to as Loschmidt’s number for this reason. In current English-language scientific literature, Loschmidt’s number is usually taken to mean the number of gas molecules in one cubic centimetre at 0° C and one atmosphere pressure (2.687×10^{19} molecules per cubic centimetre).

There are other ways by which molecular sizes and Avogadro’s number could have been estimated, such as from the spreading of a surface oil film on water or from the surface tension and the energy of evaporation of a liquid, but they will not be discussed here.

The foregoing picture of a gas as a collection of molecules

dominated by binary molecular collisions is in reality only a limited view. Two limitations of the model are briefly discussed below.

FREE-MOLECULE GAS

The mean free path in a gas may easily be increased by decreasing the pressure. If the pressure is halved, the mean free path doubles in length. Thus, at low enough pressures the mean free path can become sufficiently large that collisions of the gas molecules with surfaces become more important than collisions with other gas molecules. In such a case, the molecules can be envisioned as moving freely through space until they encounter some solid surface; hence, they are termed free-molecule gases. Such gases are sometimes called Knudsen gases, after the Danish physicist Martin Knudsen, who studied them experimentally. Many of their properties are strikingly different from those of ordinary gases (also known as continuum gases). A radiometer is a four-vaned mill that depends essentially on free-molecule effects. A temperature difference in the free-molecule gas causes a thermomolecular pressure difference that drives the vanes. The radiometer will stop spinning if enough air leaks into its glass envelope. (It will also stop spinning if all the air is removed from the envelope.) The flight of objects at high altitudes, where the mean free path is very long, is also subject to free-molecule effects. Such effects can even occur at ordinary pressures if a significant physical dimension becomes small enough. Important examples are found in many chemical process industries, where reactions are forced by catalysts to proceed at reasonable speeds. Many of these catalysts are porous materials whose pore sizes are smaller than molecular mean free paths. The speed of the desired chemical reaction may be controlled by how fast the reactant gases diffuse into the porous catalyst and by how fast the product gases can diffuse out so more reactants can enter the pores.

There is a large transition region between free-molecule behaviour and continuum behaviour, where both molecule-molecule and molecule-surface collisions are significant. This region is rather difficult to describe theoretically and remains an active field of research.

CONTINUITY OF GASEOUS AND LIQUID STATES

It may be somewhat surprising to learn that there is no fundamental distinction between a gas and a liquid. It was noted above that a gas occupies a volume about 1,600 times greater than that of an equal weight of liquid. The question arises as to the behaviour of a gas that has been compressed to 1/1,600 of its volume by application of sufficiently high pressure. If this compression is carried out above a specific temperature called the critical temperature, which is different for each gas, no phase change occurs, and the resulting substance is a gas that is just as dense as a liquid. If the compression is carried out at a fixed temperature below the critical temperature, an astonishing phenomenon occurs—at a particular pressure liquid suddenly forms. Attempts to compress the gas further simply increase the amount of liquid present and decrease the amount of gas, with the pressure remaining constant until all the gas has been converted to liquid. The applied pressure must subsequently rise a great deal to reduce the volume further, since liquids are much less compressible than gases.

Condensation

The abrupt condensation of a gas to a liquid usually does not seem astonishing because it is so commonplace—nearly everyone has boiled water, for example, which is the reverse process. From the standpoint of the kinetic-molecular theory of gases, however, it is something of a mystery. Why does it occur so abruptly and only at temperatures below a critical temperature? Equations have been written down that describe condensation, but an explanation is still lacking in the sense that no one has been able to show that it must occur, given only the forces between the molecules and the fact that their motion is described by ordinary mechanics. Condensation, which is an example of a first-order phase transition, remains one of the outstanding unsolved problems of statistical physics.

The critical temperature marks the separation between

an abrupt change and a continuous change. Other peculiar phenomena occur near the critical temperature. The densities of the coexisting liquid and gas (which is usually called a vapour in this case) become closer as the critical temperature is approached from below, and at the critical temperature they are identical. There is a unique point for every fluid, called the critical point. It is described by a critical temperature, a critical volume, and a critical pressure, at which liquid and vapour become identical. Above that temperature there is no distinction between gas and liquid; there is only a single fluid. Moreover, it is possible to pass continuously from an apparently definite gas or vapour to an apparently definite liquid with no abrupt condensation occurring. This can be accomplished by heating the vapour above the critical temperature while keeping the volume constant, then compressing it to a high density characteristic of a liquid, and finally cooling it at constant volume to its original temperature, where it is now clearly a liquid.

In short, gases and liquids are just the extreme stages of a fluid, with no fundamental distinction between the two. For this reason, an arbitrary decision has been made for the present discussion to define what is meant by the gaseous state. The definition will be based on the number density (*i.e.*, molecules per unit volume): the number density of the fluid must be low enough that only collisions between two molecules at a time need to be considered. More specifically, the mean free path must be much larger than the molecular diameter. Such a fluid shall be termed a dilute gas.

A few brief historical remarks are in order before leaving the subject of the continuity of the gaseous and liquid states. The first extensive experimental study that clearly demonstrated the phenomena involved was performed on carbon dioxide, CO₂. (Carbon dioxide, whose solid form is called dry ice, has a critical temperature of 31° C.) The experiment was conducted by Thomas Andrews at what is now the Queen's University of Belfast in Northern Ireland, and its results were summarized in 1869 in a Bakerian lecture to the Royal Society of London entitled "On the Continuity of the Gaseous and Liquid States of Matter." In 1873 a Dutch thesis was presented to the University of Leiden by Johannes D. van der Waals with virtually the same title (but in Dutch) as Andrews' lecture. In his study van der Waals used some ingenious approximations to obtain a simple equation relating the pressure, temperature, and molar volume of a fluid, based on a model that considered molecules as hard spheres with weak long-range attractive forces between them. This equation can be used to locate the critical point of a system, and it is also consistent with the occurrence of condensation when supplemented with a thermodynamic condition. This is possibly one of the most-quoted but little-read theses in science. Nevertheless, van der Waals started a scientific trend that continues to the present. His pressure-volume-temperature relation, called an equation of state, is the standard equation of state for real gases in physical chemistry, and at least one new equation of state is proposed every year in an attempt to improve on its quantitative accuracy (which is not very good). It furnished the impetus for the development of theories of liquids and of solutions. The equation is compatible with a unifying idea called the principle of corresponding states. This principle states that, if the pressure (p), volume (V), and temperature (T) of a gas are replaced, respectively, with the corresponding reduced variables—*i.e.*, the pressure divided by the critical pressure (p/p_c), the volume divided by the critical volume (V/V_c), and the temperature divided by the critical temperature (T/T_c)—all gases will behave in essentially the same manner.

The van der Waals equation

The critical point has itself proved to be a rich and deep subject. The gas-liquid critical point turns out to be only one of many types of critical points, including those of a magnetic variety, with the common feature that long-range correlations develop regardless of the molecular details of the system. That is, any small part of a system near its critical point seems to "know" what quite distant parts are doing. The mathematical description of the behaviour of a system near its critical point also becomes rather unusual.

Behaviour and properties

Statistical methods

The enormous number of molecules in even a small volume of a dilute gas produces not complication, as might be expected, but rather simplification. The reason is that ordinarily only statistical averages are observed in the study of the behaviour and properties of gases, and statistical methods are quite accurate when large numbers are involved. Compared to the numbers of molecules involved, there are only a few properties of gases that warrant attention here, namely, pressure, density, temperature, internal energy, viscosity, heat conductivity, and diffusivity. (More subtle properties can be brought into view by the application of electric and magnetic fields, but they are of minor interest.)

It is a remarkable fact that these properties are not independent. If two are known, the rest can be determined from them. That is to say, for a given gas, the specification of only two properties—usually chosen to be temperature and density or temperature and pressure—fixes all the others. Thus, if the temperature and density of carbon dioxide are specified, the gas can have only one possible pressure, one internal energy, one viscosity, and so on. In order to determine the values of these other properties, they must either be measured or calculated from the known properties of the molecules themselves. Such calculations are the ultimate goal of statistical mechanics and kinetic theory, and dilute gases constitute the case for which the most progress toward that goal has been made.

In discussing the behaviour of gases, it is useful to separate the equilibrium properties and the nonequilibrium transport properties. By definition, a system in equilibrium can undergo no net change unless some external action is performed on it (*e.g.*, pushing in a piston or adding heat). Its behaviour is steady with time, and no changes appear to be occurring, even though the molecules are in ceaseless motion. In contrast, the nonequilibrium properties describe how a system responds to some external action, such as the imposition of a temperature or pressure difference. Equilibrium behaviour is much easier to analyze, because any change that occurs on the molecular level must be compensated by some other change or changes on the molecular level in order for the system to remain in equilibrium.

EQUILIBRIUM PROPERTIES

Among the most obvious properties of a dilute gas, other than its low density compared with liquids and solids, are its great elasticity or compressibility and its large volume expansion on heating. These properties are nearly the same for all dilute gases, and virtually all such gases can be described quite accurately by the following universal equation of state:

$$pv = RT. \quad (15)$$

This expression is called the ideal, or perfect, gas equation of state, since all real gases show small deviations from it, although these deviations become less significant as the density is decreased. Here p is the pressure, v is the volume per mole, or molar volume, R is the universal gas constant, and T is the absolute thermodynamic temperature. To a rough degree, the expression is accurate within a few percent if the volume is more than 10 times the critical volume; the accuracy improves as the volume increases. The expression eventually fails at both high and low temperatures, owing to ionization at high temperatures and to condensation to a liquid or solid at low temperatures.

The ideal gas equation of state is an amalgamation of three ideal gas laws that were formulated independently. The first is Boyle's law, which refers to the elastic properties of the gas; it was described by the Anglo-Irish scientist Robert Boyle in 1662 in his famous "... Experiments ... Touching the Spring of the Air ...". It states that the volume of a gas at constant temperature is inversely proportional to the pressure; *i.e.*, if the pressure on a gas is doubled, for example, its volume decreases by one-half. The second, usually called Charles's law, is concerned with the thermal expansion of the gas. It is named in honour of the French experimental physicist Jacques-

Alexandre-César Charles for the work he carried out in about 1787. The law states that the volume of a gas at constant pressure is directly proportional to the absolute temperature; *i.e.*, an increase of temperature of 1° C at room temperature causes the volume to increase by about 1 part in 300, or 0.3 percent. The third law embodied in equation (15) is based on the 1811 hypothesis of the Italian scientist Amedeo Avogadro—namely, that equal volumes of gases at the same temperature and pressure contain equal numbers of particles. The number of particles (or molecules) is proportional to the number of moles n , the constant of proportionality being Avogadro's number, N_0 . Thus, at constant temperature and pressure the volume of a gas is proportional to the number of moles. If the total volume V contains n moles of gas, then only $v = V/n$ appears in the equation of state. By measuring the quantity of gas in moles rather than grams, the constant R is made universal; if mass were measured in grams (and hence v in volume per gram), then R would have a different value for each gas.

The ideal gas law is easily extended to mixtures by letting n represent the total number of moles of all species present in volume V . That is, if there are n_1 moles of species 1, n_2 moles of species 2, etc., in the mixture, then $n = n_1 + n_2 + \dots$ and $v = V/n$ as before. This result can also be rewritten and reinterpreted in terms of the partial pressures of the different species, such that $p_1 = n_1RT/V$ is the partial pressure of species 1 and so on. The total pressure is then given as $p = p_1 + p_2 + \dots$. This rule is known as Dalton's law of partial pressures in honour of the British chemist and physicist John Dalton, who formulated it about 1801.

A brief aside on units and temperature scales is in order. The (metric) unit of pressure in the scientific international system of units (known as the SI system) is newton per square metre (N/m^2), where one newton (N) is the force that gives a mass of one kilogram an acceleration of 1 m/s^2 . The unit N/m^2 is given the name pascal (Pa), where one standard atmosphere is exactly 101,325 Pa (approximately 14.7 pounds per square inch). The unit of volume in the SI system is the cubic metre ($1 \text{ m}^3 = 10^6 \text{ cm}^3$), and the unit of temperature is the kelvin (K). The Kelvin thermodynamic temperature scale is defined through the laws of thermodynamics so as to be absolute or universal, in the sense that its definition does not depend on the specific properties of any particular kind of matter. Its numerical values, however, are assigned by defining the triple point of water—*i.e.*, the unique temperature at which ice, liquid water, and water vapour are all in equilibrium—to be exactly 273.16 K. The freezing point of water under one atmosphere of air then turns out to be (by measurement) 273.1500 K. The freezing point is 0° on the Celsius scale (or 32° on the Fahrenheit scale), by definition. The precise thermodynamic definition of the Kelvin scale and the rather peculiar number chosen to define its numerical values (*i.e.*, 273.16) are historical choices made so that the ideal gas equation of state will have the simple mathematical form given by the right-hand side of equation (15).

Kelvin scale

The gas constant R is determined by measurement. The best value so far obtained is that of the U.S. National Institute of Standards and Technology—namely, 8.314471 J/mol · K. Here the unit J is one of work or energy, one joule (J) being equal to one newton-metre.

Once the equation of state is known for an ideal gas, only its internal energy, E , needs to be determined in order for all other equilibrium properties to be deducible from the laws of thermodynamics. That is to say, if the equation of state and the internal energy of a fluid are known, then all the other thermodynamic properties (*e.g.*, enthalpy, entropy, and free energy) are fixed by the condition that it must be impossible to construct perpetual motion machines from the fluid. Proofs of such statements are usually rather subtle and involved and constitute a large part of the subject of thermodynamics, but conclusions based on thermodynamic principles are among the most reliable results of science.

A thermodynamic result of relevance here is that the ideal gas equation of state requires that the internal energy depend on temperature alone, not on pressure or density.

Ideal gas equation of state

The actual relationship between E and T must be measured or calculated from known molecular properties by means of statistical mechanics. The internal energy is not directly measurable, but its behaviour can be determined from measurements of the molar heat capacity (*i.e.*, the specific heat) of the gas. The molar heat capacity is the amount of energy required to raise the temperature of one mole of a substance by one degree; its units in the SI system are $\text{J/mol} \cdot \text{K}$. A system with many kinds of motion on a molecular scale absorbs more energy than one with only a few kinds of motion. The interpretation of the temperature dependence of E is particularly simple for dilute gases, as is shown in the discussion of the kinetic theory of gases below. The following highlights only the major aspects.

Every gas molecule moves in three-dimensional space, and this translational motion contributes $(3/2)RT$ (per mole) to the internal energy E . For monatomic gases, such as helium, neon, argon, krypton, and xenon, this is the sole energy contribution. Gases that contain two or more atoms per molecule also contribute additional terms because of their internal motions:

$$E \text{ (per mole)} = \frac{3}{2}RT + E_{int} \quad (16)$$

where E_{int} may include contributions from molecular rotations and internal vibrations and occasionally from internal electronic excitations. Some of these internal motions may not contribute at ordinary temperatures because of special conditions imposed by quantum mechanics, however, so that the temperature dependence of E_{int} can be rather complex.

The extension to gas mixtures is straightforward—the total internal energy E (per mole) is the weighted sum of the internal energies of each of the species: $nE = n_1E_1 + n_2E_2 + \dots$, where $n = n_1 + n_2 + \dots$.

It is the task of the kinetic theory of gases to account for these results concerning the equation of state and the internal energy of dilute gases.

TRANSPORT PROPERTIES

The following is a summary of the three main transport properties: viscosity, heat conductivity, and diffusivity. These properties correspond to the transfer of momentum, energy, and matter, respectively.

Viscosity. All ordinary fluids exhibit viscosity, which is a type of internal friction. A continuous application of force is needed to keep a fluid flowing, just as a continuous force is needed to keep a solid body moving in the presence of friction. Consider the case of a fluid slowly flowing through a long capillary tube. A pressure difference of Δp must be maintained across the ends to keep the fluid flowing, and the resulting flow rate is proportional to Δp . The rate is inversely proportional to the viscosity (η) since the friction that opposes the flow increases as η increases. It also depends on the geometry of the tube, but this effect will not be considered here. The SI units of η are $\text{N} \cdot \text{s}/\text{m}^2$ or $\text{Pa} \cdot \text{s}$. An older unit of the centimetre-gram-second version of the metric system that is still often used is the poise ($1 \text{ Pa} \cdot \text{s} = 10 \text{ poise}$). At 20°C the viscosity of water is $1.0 \times 10^{-3} \text{ Pa} \cdot \text{s}$ and that of air is $1.8 \times 10^{-5} \text{ Pa} \cdot \text{s}$. To a rough approximation, liquids are about 100 times more viscous than gases.

There are three important properties of the viscosity of dilute gases that seem to defy common sense. All can be explained, however, by the kinetic theory (see below *Kinetic theory of gases*). The first property is the lack of a dependence on pressure or density. Intuition suggests that gas viscosity should increase with increasing density, inasmuch as liquids are much more viscous than gases, but gas viscosity is actually independent of density. This result can be illustrated by a pendulum swinging on a solid support. It eventually slows down owing to the viscous friction of the air. If a bell jar is placed over the pendulum and half the air is pumped out, the air remaining in the jar damps the pendulum just as fast as a full jar of air would have done. Robert Boyle noted this peculiar phenomenon in 1660, but his results were largely either ignored or forgotten. The Scottish chemist Thomas Graham studied

the flow of gases through long capillaries, which he called transpiration, in 1846 and 1849, but it was not until 1877 that the German physicist O.E. Meyer pointed out that Graham's measurements had shown the independence of viscosity on density. Prior to Meyer's investigations, the kinetic theory had suggested the result, so he was looking for experimental proof to support the prediction. When James Clerk Maxwell discovered (in 1865) that his kinetic theory suggested this result, he found it difficult to believe and attempted to check it experimentally. He designed an oscillating disk apparatus (which is still much copied) to verify the prediction.

The second unusual property of viscosity is its relationship with temperature. One might expect the viscosity of a fluid to increase as the temperature is lowered, as suggested by the phrase "as slow as molasses in January." The viscosity of a dilute gas behaves in exactly the opposite way: the viscosity increases as the temperature is raised. The rate of increase varies approximately as T^s , where s is between $1/2$ and 1, and depends on the particular gas. This behaviour was clearly established in 1849 by Graham.

The third property pertains to the viscosity of mixtures. A viscous syrup, for example, can be made less so by the addition of a liquid with a lower viscosity, such as water. By analogy, one would expect that a mixture of carbon dioxide, which is fairly viscous, with a gas like hydrogen, which is much less viscous, would have a viscosity intermediate to that of carbon dioxide and hydrogen. Surprisingly, the viscosity of the mixture is even greater than that of carbon dioxide. This phenomenon was also observed by Graham in 1849.

Finally, there is no obvious correlation of gas viscosity with molecular weight. Heavy gases are often more viscous than light gases, but there are many exceptions, and no simple pattern is apparent.

Heat conduction. If a temperature difference is maintained across a fluid, a flow of energy through the fluid will result. The energy flow is proportional to the temperature difference according to Fourier's law, where the constant of proportionality (aside from the geometric factors of the apparatus) is called the heat conductivity or thermal conductivity of the fluid, λ . Mechanisms other than conduction can transport energy, in particular convection and radiation; here it is assumed that these can be eliminated or adjusted for. The SI units for λ are $\text{J}/\text{m} \cdot \text{s} \cdot \text{K}$ or watt per metre degree ($\text{W}/\text{m} \cdot \text{K}$), but sometimes calories are used for the energy term instead of joules (one calorie = 4.184 J). At 20°C the thermal conductivity of water is $0.60 \text{ W}/\text{m} \cdot \text{K}$, and that of many organic liquids is roughly only one-third as large. The thermal conductivity of air at 20°C is only about $2.5 \times 10^{-2} \text{ W}/\text{m} \cdot \text{K}$. To a rough approximation, liquids conduct heat about 10 times better than do gases.

The properties of the thermal conductivity of dilute gases parallel those of viscosity in some respects. The most striking is the lack of dependence on pressure or density. Based on this fact, there seems to be no advantage to pumping out the inner chambers of thermos bottles. As far as conduction is concerned, it does not provide any benefits until practically all the air has been removed and free-molecule conduction is occurring. Convection, however, does depend on density, so some degree of insulation is provided by pumping out only some of the air.

The thermal conductivity of a dilute gas increases with increasing temperature, much like its viscosity. In this case, such behaviour does not seem particularly odd, probably because most people do not have a preconceived idea of how thermal conductivity should behave, unlike the situation with viscosity.

There are some differences in the behaviour of thermal conductivity and viscosity; one of the most notable has to do with mixtures. At first glance the thermal conductivity of a gaseous mixture seems to be as expected, since it falls between the conductivities of its components, but a closer look reveals an odd regularity. The conductivity of the mixture is always less than an average based on the number of moles (or molecules) of each component in the mixture. This appears to be related to the different effect that molecular weight has on thermal conductivity and

Contributions to internal energy

Temperature-viscosity relationship

Conductivity of mixtures

viscosity. Light gases are usually better conductors than are heavy gases, whereas heavy gases are often (but not always) more viscous than are light gases. There also seems to be some correlation between molar heat capacity and thermal conductivity. The foregoing properties of thermal conductivity pose more puzzles that the kinetic theory of gases must address.

Diffusion. Diffusion in dilute gases is in some ways more complex, or at least more subtle, than either viscosity or thermal conductivity. First, a mixture is necessarily involved, inasmuch as a gas diffusing through itself makes no sense physically unless the molecules are in some way distinguishable from one another. Second, diffusion measurements are rather sensitive to the details of the experimental conditions. This sensitivity can be illustrated by the following considerations.

Light molecules have higher average speeds than do heavy molecules at the same temperature. This result follows from kinetic theory, as explained below, but it can also be seen by noting that the speed of sound is greater in a light gas than in a heavy gas. This is the basis of the well-known demonstration that breathing helium causes one to speak with a high-pitched voice. If a light and a heavy gas are interdiffusing, the light molecules should move into the heavy-gas region faster than the heavy molecules move into the light-gas region, thereby causing the pressure to rise in the heavy-gas region. If the diffusion takes place in a closed vessel, the pressure difference drives the heavy gas into the light-gas region at a faster rate than it would otherwise diffuse, and a steady state is quickly reached in which the number of heavy molecules traveling in one direction equals, on the average, the number of light molecules traveling in the opposite direction. This method, called equimolar countercurrent diffusion, is the usual manner in which gaseous diffusion measurements are now carried out.

The steady-state pressure difference that develops is almost unmeasurably small unless the diffusion occurs through a fine capillary or a fine-grained porous material. Nevertheless, experimenters have been able to devise clever schemes either to measure it or to prevent its development. The first to do the latter was Graham in 1831; he kept the pressure uniform by allowing the gas mixture to flow. The results of this work now appear in elementary textbooks as Graham's law of diffusion. Most of these accounts are incorrect or incomplete or both, owing to the fact that the writers confuse the uniform-pressure experiment either with the equal countercurrent experiment or with the phenomenon of effusion (described below in the section *Kinetic theory of gases*). Graham also performed equal countercurrent experiments in 1863, using a long closed-tube apparatus he devised. This sort of apparatus is now usually called a Loschmidt diffusion tube after Loschmidt, who used a modified version of the tube in 1870 to make a series of accurate diffusion measurements on a number of gas pairs.

A quantitative description of diffusion follows. A composition difference in a two-component gas mixture causes a relative flow of the components that tends to make the composition uniform. The flow of one component is proportional to its concentration difference, and in an equal countercurrent experiment this is balanced by an equal and opposite flow of the other component. The constant of proportionality is the same for both components and is called the diffusion coefficient, D_{12} , for that gas pair. This relationship between the flow rate and the concentration difference is called Fick's law of diffusion. The SI units for the diffusion coefficient are square metres per second (m^2/s). Diffusion, even in gases, is an extremely slow process, as was pointed out above in estimating molecular sizes and collision rates. Gaseous diffusion coefficients at one atmosphere pressure and ordinary temperatures lie largely in the range of 10^{-5} to 10^{-4} m^2/s , but diffusion coefficients for liquids and solutions lie in the range of only 10^{-10} to 10^{-9} m^2/s . To a rough approximation, gases diffuse about 100,000 times faster than do liquids.

Diffusion coefficients are inversely proportional to total pressure or total molar density and are therefore reported by convention at a standard pressure of one atmosphere.

Doubling the pressure of a diffusing mixture halves the diffusion coefficient, but the actual rate of diffusion remains unchanged. This seemingly paradoxical result occurs because doubling the pressure also doubles the concentration, according to the ideal gas equation of state, and hence doubles the concentration difference, which is the driving force for diffusion. The two effects exactly compensate.

Diffusion coefficients increase with increasing temperature at a rate that depends on whether the pressure or the total molar density is held constant as the temperature is changed. If the rate increases as T^s at constant molar density (where s usually lies between $1/2$ and 1), then it will increase as T^{1+s} at constant pressure, according to the ideal gas equation of state.

Perhaps the most surprising property of gaseous diffusion coefficients is that they are virtually independent of the mixture's composition, varying by at most a few percent over the whole composition range, even for very dissimilar gases. A trace of hydrogen, for example, diffuses through carbon dioxide at virtually the same rate that a trace of carbon dioxide diffuses through hydrogen. Liquid mixtures do not behave this way, and liquid diffusion coefficients may vary by as much as a factor of 10 from one end of the composition range to the other. The lack of composition dependence of gaseous diffusion coefficients is one of the odder properties to be explained by kinetic theory.

Thermal diffusion. If a temperature difference is applied to a uniform mixture of two gases, the mixture will partially separate into its components, with the heavier, larger molecules usually (but not invariably) concentrating at the lower temperature. This behaviour was predicted theoretically before it was observed experimentally, but a rather elaborate explanation was required because simple theory suggests no such phenomenon. It was predicted in 1911–12 by David Enskog in Sweden and independently in 1917 by Sydney Chapman in England, but the validity of their theoretical results was questioned until Chapman (who was an applied mathematician) enlisted the aid of the chemist F.W. Dootson to verify it experimentally.

Thermal diffusion can be used to separate isotopes. The amount of separation for any reasonable temperature difference is quite small for isotopes, but the effect can be amplified by combining it with slow thermal convection in a columnar arrangement devised in 1938 by Klaus Clusius and Gerhard Dickel in Germany. While the apparatus is quite simple, the theory of its operation is not: a long cylinder with a diameter of several centimetres is mounted vertically with an electrically heated hot wire along its central axis. The thermal diffusion occurs horizontally between the hot wire and the cold wall of the cylinder, and the convection takes place vertically to bring new gas regions into contact.

There is also an effect that is the inverse of thermal diffusion, called the diffusion thermoeffect, in which an imposed concentration difference causes a temperature difference to develop. That is, a diffusing gas mixture develops small temperature differences, on the order of 1°C , which die out as the composition approaches uniformity. The transport coefficient describing the diffusion thermoeffect must be equal to the coefficient describing thermal diffusion, according to the reciprocal relations central to the thermodynamics of irreversible processes.

The diffusion thermoeffect

Kinetic theory of gases

The aim of kinetic theory is to account for the properties of gases in terms of the forces between the molecules, assuming that their motions are described by the laws of mechanics (usually classical Newtonian mechanics, although quantum mechanics is needed in some cases). The present discussion focuses on dilute ideal gases, in which molecular collisions of at most two bodies are of primary importance. Only the simplest theories are treated here in order to avoid obscuring the fundamental physics with complex mathematics.

IDEAL GAS

The ideal gas equation of state can be deduced by calculating the pressure as caused by molecular impacts on a

container wall. The internal energy and Dalton's law of partial pressures also emerge from this calculation, along with some free-molecule phenomena. The calculation is significant because it is basically the same one used to explain all dilute-gas phenomena.

Pressure. Newton's second law of motion can be stated in not-so-familiar form as impulse equals change in momentum, where impulse is force multiplied by the time during which it acts. A molecule experiences a change in momentum when it collides with a container wall; during the collision an impulse is imparted by the wall to the molecule that is equal and opposite to the impulse imparted by the molecule to the wall. This is required by Newton's third law. The sum of the impulses imparted by all the molecules to the wall is, in effect, the pressure. Consider a system of molecules of mass m traveling with a velocity v in an enclosed container. In order to arrive at an expression for the pressure, a calculation will be made of the impulse imparted to one of the walls by a single impact, followed by a calculation of how many impacts occur on that wall during a time t . Although the molecules are moving in all directions, only those with a component of velocity toward the wall can collide with it; call this component v_z , where z represents the direction directly toward the wall. Not all molecules have the same v_z , of course; perhaps only N_z out of a total of N molecules do. To find the total pressure, the contributions from molecules with all different values of v_z must be summed. A molecule approaches the wall with an initial momentum mv_z , and after impact it moves away from the wall with an equal momentum in the opposite direction, $-mv_z$. Thus, the total change in momentum is $mv_z - (-mv_z) = 2mv_z$, which is equal to the total impulse imparted to the wall.

The number of impacts on a small area A of the wall in time t is equal to the number of molecules that reach the wall in time t . Since the molecules are traveling at speed v_z , only those within a distance $v_z t$ and moving toward the wall will reach it in that time. Thus, the molecules that are traveling toward the wall and are within a volume $Av_z t$ will strike the area A of the wall in time t . On the average, half of the molecules in this volume will be moving toward the wall. If N_z molecules with speed component v_z are present in the total volume V , then $(1/2)(N_z/V)(A)(v_z t)$ molecules in the collision volume will hit, and each one contributes an impulse of $2mv_z$. The total impulse in time t is therefore $(1/2)(N_z/V)(A)(v_z t)(2mv_z) = (N_z/V)(mv_z^2)(At)$, which is equal to Ft , where F is the force on the wall due to the impacts. Equating these two expressions, the time factor t cancels out. Since pressure is defined as the force per unit area (F/A), it follows that the contribution to the pressure from the molecules with speed v_z is thus $(N_z/V)mv_z^2$. Because there are different values of v_z^2 for different molecules, the average value, denoted \bar{v}_z^2 , is used to take into account the contributions from all the molecules. The pressure is thus given as $p = (N/V)mv_z^2$.

Since the molecules are in random motion, this result is independent of the choice of axis. For any choice of (x , y , z) axes, the magnitude of the velocity is $v^2 = v_x^2 + v_y^2 + v_z^2$ (which is just the Pythagorean theorem in three dimensions), and taking the average gives $v^2 = v_x^2 + v_y^2 + v_z^2$. The gas is in equilibrium, so it must appear the same in any direction, and the average velocities are therefore the same in all directions—*i.e.*, $v_x^2 = v_y^2 = v_z^2$; thus $v^2 = 3v_z^2$. When the value $(1/3)v^2$ is substituted for v_z^2 in the expression for pressure, the following equation is obtained:

$$p = \frac{1}{3} \frac{N}{V} m \bar{v}^2, \quad \text{or} \quad pV = \frac{1}{3} N m \bar{v}^2. \quad (17)$$

To rewrite this in molar units, N is set equal to nN_0 —*i.e.*, the product of the number of moles n and Avogadro's number N_0 —to give

$$pV = \frac{1}{3} n M \bar{v}^2, \quad (18)$$

where $M = N_0 m$ is the molecular weight of the gas and v is the volume per mole (V/n). Since the ideal gas equation of state relates pressure, molar volume, and temperature as $pV = RT$, the temperature T must be related to the average kinetic energy of the molecules as

$$\frac{1}{2} M \bar{v}^2 = \frac{3}{2} RT. \quad (19)$$

This expression is often written in molecular (rather than molar) terms as $(1/2)[m\bar{v}^2] = (3/2)kT$, where $k = R/N_0$ is called Boltzmann's constant. If the gas is a mixture, the foregoing calculation shows that the impacts of the different species are simply added separately, and Dalton's law of partial pressures follows directly.

The energy law given as equation (16) also follows from equation (19): the kinetic energy of translational motion per mole is $(3/2)RT$. Any energy residing in the internal motions of the individual molecules is simply carried separately without contributing to the pressure.

Average molecular speeds can be calculated from the results of kinetic theory in terms of the so-called root-mean-square speed v_{rms} . The v_{rms} is the square root of the average of the squares of the speeds of the molecules: $(v^2)^{1/2}$. From equation (19) the v_{rms} is $(3RT/M)^{1/2}$. At 20° C the value for air ($M = 29$) is 502 m/s, a result very close to the rough estimate of 5×10^2 m/s given above.

Molecule-molecule collisions were not considered in the calculation of the expression for pressure even though many such collisions occur. Such collisions could be ignored because they are elastic; *i.e.*, linear momentum is conserved in the collision, provided that no external forces act. Two molecules therefore continue to carry the same momentum to the wall even if they collide with one another before striking it. The ideal gas equation of state remains valid as the density is decreased, even holding for a free-molecule gas. The equation eventually fails as the density is increased, however, because other molecules exert forces and change the rate of collisions with the walls.

It was not until the mid- to late 19th century that kinetic theory was successfully applied to such calculations as gas pressure. Such notable scientists as Sir Isaac Newton and John Dalton had believed that gas pressure was caused by repulsions between molecules that pushed them against the container walls. For many reasons, the kinetic theory had overshadowed such static theories (and others such as vortex theories) by about 1860. It was not until 1875, however, that Maxwell actually proved that a static theory was in conflict with experiment.

Effusion. Consider the system described above in the calculation of gas pressure, but with the area A in the container wall replaced with a small hole. The number of molecules that escape through the hole in time t is equal to $(1/2)(N/V)\bar{v}_z(At)$. In this case, collisions between molecules are significant, and the result holds only for tiny holes in very thin walls (as compared to the mean free path), so that a molecule that approaches near the hole will get through without colliding with another molecule and being deflected away. The relationship between v_z and the average speed \bar{v} is rather straightforward: $\bar{v}_z = (1/2)\bar{v}$.

If the rates for two different gases effusing through the same hole are compared, starting with the same gas density each time, it is found that much more light gas escapes than heavy gas and that more gas escapes at a high temperature than at a low temperature, other things being equal. In particular,

$$\frac{\text{effusion rate of gas 1}}{\text{effusion rate of gas 2}} = \frac{\bar{v}_1}{\bar{v}_2} = \left(\frac{m_2}{m_1}\right)^{1/2} \left(\frac{T_1}{T_2}\right)^{1/2}. \quad (20)$$

The last step follows from the energy formula, $(1/2)m\bar{v}^2 = (3/2)kT$, where $(v^2)^{1/2}$ is approximated to be v , even though v^2 and $(\bar{v})^2$ actually differ by a numerical factor near unity (namely, $3\pi/8$). This result was discovered experimentally in 1846 by Graham for the case of constant temperature and is known as Graham's law of effusion. It can be used to measure molecular weights, to measure the vapour pressure of a material with a low vapour pressure, or to calculate the rate of evaporation of molecules from a liquid or solid surface.

Thermal transpiration. Suppose that two containers of the same gas but at different temperatures are connected by a tiny hole and that the gas is brought to a steady state. If the hole is small enough and the gas density is low enough that only effusion occurs, the equilibrium pressure will be greater on the high-temperature side. But, if the

Molecule-wall collisions

Early static gas theory

Average velocity

Graham's law of effusion

initial pressures on both sides are equal, gas will flow from the low-temperature side to the high-temperature side to cause the high-temperature pressure to increase. The latter situation is called thermal transpiration, and the steady-state result is called the thermomolecular pressure difference. These results follow simply from the effusion formula if the ideal gas law is used to replace N/V with p/T :

$$\frac{\text{effusion rate from container 1}}{\text{effusion rate from container 2}} = \frac{p_1}{p_2} \left(\frac{T_2}{T_1} \right)^{1/2}. \quad (21)$$

When a steady state is reached, the effusion rates are equal, and thus

$$\frac{p_1}{p_2} = \left(\frac{T_1}{T_2} \right)^{1/2}. \quad (22)$$

This phenomenon was first investigated experimentally by Osborne Reynolds in 1879 in Manchester, Eng. Errors can result if a gas pressure is measured in a vessel at very low or very high temperature by connecting it via a fine tube to a manometer at room temperature. A continuous circulation of gas can be produced by connecting the two containers with another tube whose diameter is large compared with the mean free path. The pressure difference drives gas through this tube by viscous flow. A continuous engine based on this circulating flow unfortunately has a low efficiency.

Viscosity. The kinetic-theory explanation of viscosity can be simplified by examining it in qualitative terms. Viscosity is caused by the transfer of momentum between two planes sliding parallel to one another but at different rates, and this momentum is transferred by molecules moving between the planes. Molecules from the faster plane move to the slower plane and tend to speed it up, while molecules from the slower plane travel to the faster plane and tend to slow it down. This is the mechanism by which one plane experiences the drag of the other. A simple analogy is two mail trains passing each other, with workers throwing mailbags between the trains. Every time a mailbag from the fast-moving train lands on the slow one, it imparts its momentum to the slow train, speeding it up a little; likewise each mailbag from the slow train that lands on the fast one slows it down a bit.

If the trains are too far apart, the mailbags cannot be passed between them. Similarly, the planes of a gas must be only about a mean free path apart in order for molecules to pass between them without being deflected by collisions. If one uses this approach, a simple calculation can be carried out, much as in the case of the gas pressure, with the result that

$$\eta = a \frac{N}{V} \bar{v} l m, \quad (23)$$

where a is a numerical constant of order unity, the term $(N/V)\bar{v}l$ is a measure of the number of molecules contained in a small counting cylinder, and the mass m is a measure of the momentum carried between the sliding planes. The cross-sectional area of the counting cylinder and the relative speed of the sliding planes do not appear in the equation because they cancel one another when the drag force is divided by the area and speed of the planes in order to find η .

It can now be seen why η is independent of gas density or pressure. The term (N/V) in equation (23) is the number of carriers of momentum, but l measures the number of collisions that interfere with these carriers and is inversely proportional to (N/V) . The two effects exactly cancel each other. Viscosity increases with temperature because the average velocity \bar{v} does; that is, momentum is carried more quickly when the molecules move faster. Although \bar{v} increases as $T^{1/2}$, η increases somewhat faster because the mean free path also increases with temperature, since it is harder to deflect a fast molecule than a slow one. This feature depends explicitly on the forces between the molecules and is difficult to calculate accurately, as is the value of the constant a , which turns out to be close to $1/2$.

The behaviour of the viscosity of a mixture can also be explained by the foregoing calculation. In a mixture of a light gas and a viscous heavy gas, both types of molecules have the same average energy; however, most of the mo-

mentum is carried by the heavy molecules, which are therefore the main contributors to the viscosity. The light molecules are rather ineffective in deflecting the heavy molecules, so that the latter continue to carry virtually as much momentum as they would in the absence of light molecules. The addition of a light gas to a heavy gas therefore does not reduce the viscosity substantially and may in fact increase it because of the small extra momentum carried by the light molecules. The viscosity will eventually decrease when there are only a few heavy molecules remaining in a large sea of light molecules.

The main dependence of η on the molecular mass is through the product $\bar{v}m$ in equation (23), which varies as $m^{1/2}$ since \bar{v} varies as $1/m^{1/2}$. Owing to this effect, heavy gases tend to be more viscous than light gases, but this tendency is compensated for to some degree by the behaviour of l , which tends to be smaller for heavy molecules because they are usually larger than light molecules and therefore more likely to collide. The often confusing connection between viscosity and molecular weight can thus be accounted for by equation (23).

Finally, in a free-molecule gas there are no collisions with other molecules to impede the transport of momentum, and the viscosity thus increases linearly with pressure or density until the number of collisions becomes great enough so that the viscosity assumes the constant value given by equation (23). The nonideal behaviour of the gas that accompanies further increases in density eventually leads to an increase in viscosity, and the viscosity of an extremely dense gas becomes much like that of a liquid.

Thermal conductivity. The kinetic-theory explanation of heat conduction is similar to that for viscosity, but in this case the molecules carry net energy from a region of higher energy (*i.e.*, temperature) to one of lower energy (temperature). Internal molecular motions must be accounted for because, though they do not transport momentum, they do transport energy. Monatomic gases, which carry only their kinetic energy of translational motion, are the simplest case. The resulting expression for thermal conductivity is

$$\lambda = a' \frac{N}{V} \bar{v} l \left(\frac{3}{2} k \right), \quad (24)$$

which has the same basic form as equation (23) for viscosity, with $(3k/2)$ replacing m . The $(3k/2)$ is the heat capacity per molecule and is the conversion factor from an energy difference to a temperature difference.

It can be shown from equation (24) that the independence of density and the increase with temperature is the same for thermal conductivity as it is for viscosity. The dependence on molecular mass is different, however, with λ varying as $1/m^{1/2}$ owing to the factor \bar{v} . Thus, light gases tend to be better conductors of heat than are heavy gases, and this tendency is usually augmented by the behaviour of l .

The behaviour of the thermal conductivity of mixtures may be qualitatively explained. Adding heavy gas to light gas reduces the thermal conductivity because the heavy molecules carry less energy and also interfere with the energy transport of the light molecules.

The similar behaviour of λ and η suggests that their ratio might provide information about the constants a and a' . The ratio of a'/a is given as

$$\frac{\lambda}{\eta} \frac{m}{(3k/2)} = \frac{a'}{a}. \quad (25)$$

Although simple theory suggests that this ratio should be about one, both experiment and more refined theory give a value close to $5/2$. This means that molecules do not "forget" their past history in every collision, but some persistence of their precollision velocities occurs. Molecules transport both energy and momentum from a somewhat greater distance than just one mean free path, but this distance is greater for energy than for momentum. This is plausible, for molecules with higher kinetic energies might be expected to have greater persistences.

Attempts to calculate the constants a and a' by tracing collision histories to find the "persistence of velocities" have not met with much success. The molecular "mem-

ory" fades slowly, too many previous collisions have to be traced, and the calculations become almost hopelessly complicated. A different theoretical approach is needed, which was finally supplied about 1916–17 independently by Enskog and Chapman. Their theory also shows that the same value of l applies to both η and λ , a fact that is not obvious in the simple theory described here.

The thermal conductivity of polyatomic molecules is accounted for by simply adding on a contribution for the energy carried by the internal molecular motions:

$$\lambda = a' \left(\frac{N}{V} \right) \bar{v} l \left(\frac{3}{2} k \right) + a'' \left(\frac{N}{V} \right) \bar{v} l c_{int}, \quad (26)$$

where c_{int} is the contribution of the internal motions to the heat capacity (per molecule) and is easily found by subtracting $(3k/2)$ from the total measured heat capacity. As might be expected, the constant a'' is only about half as large as a' .

The pressure or density dependence of λ must be similar to that of η —an initial linear increase in the free-molecule region, followed by a constant value in the dilute-gas region and finally an increase in the dense-fluid region.

Difficulties
in applying
simple
theory

Diffusion and thermal diffusion. Both of these properties present difficulties for the simple mean free path version of kinetic theory. In the case of diffusion it must be argued that collisions of the molecules of species 1 with other species 1 molecules do not inhibit the interdiffusion of species 1 and 2, and similarly for 2–2 collisions. If this is not assumed, the calculated value of the diffusion coefficient for the 1–2 gas pair, D_{12} , depends strongly on the mixture composition instead of being virtually independent of it, as is shown by experiment. The neglect of 1–1 and 2–2 collisions can be rationalized by noting that the flow of momentum is not disturbed by such like-molecule collisions owing to the conservation of momentum, but it can be contended that the argument was simply invented to make the theory agree with experiment. A more charitable view is that the experimental results demonstrate that collisions between like molecules have little effect on D_{12} . It is one of the triumphs of the accurate kinetic theory of Enskog and Chapman that this result clearly emerges.

If 1–1 and 2–2 collisions are ignored, a simple calculation gives a result much like those for η and λ :

$$D_{12} = a_{12} \bar{v}_{12} l_{12}, \quad (27)$$

where a_{12} is a numerical constant, \bar{v}_{12} is an average relative speed for 1–2 collisions given by $\bar{v}_{12}^2 = (1/2)(\bar{v}_1^2 + \bar{v}_2^2)$, and l_{12} is a mean free path for 1–2 collisions that is inversely proportional to the total molecular number density, $(N_1 + N_2)/V$. Thus, D_{12} is inversely proportional to gas density or pressure, unlike η and λ , but the concentration difference is proportional to pressure, with the two effects canceling one another, as pointed out previously. The actual transport of molecules is therefore independent of pressure. The numerical value of a_{12} , as obtained by refined calculations, is close to 3/5.

The pressure dependence of pD_{12} should be qualitatively similar to that of η and λ —an initial linear increase in the free-molecule region, a constant value in the dilute-gas region, and finally an increase in the dense-fluid region.

Thermal diffusion presents special difficulties for kinetic theory. The transport coefficients η , λ , and D_{12} are always positive regardless of the nature of the intermolecular forces that produce the collisions—the mere existence of collisions suffices to account for their important features. The transport coefficient that describes thermal diffusion, however, depends critically on the nature of the intermolecular forces and the collisions and can be positive, negative, or zero. Its dependence on composition is also rather complicated. There have been a number of attempts to explain thermal diffusion with a simple mean free path model, but none has been satisfactory. No simple physical explanation of thermal diffusion has been devised, and recourse to the accurate, but complicated, kinetic theory is necessary.

Boltzmann equation. The simple mean free path description of gas transport coefficients accounts for the major observed phenomena, but it is quantitatively unsatisfactory with respect to two major points: the values

of numerical constants such as a , a' , a'' , and a_{12} and the description of the molecular collisions that define a mean free path. Indeed, collisions remain a somewhat vague concept except when they are considered to take place between molecules modeled as hard spheres. Improvement has required a different, somewhat indirect, and more mathematical approach through a quantity called the velocity distribution function. This function describes how molecular velocities are distributed on the average: a few very slow molecules, a few very fast ones, and most near some average value—namely, $v_{rms} = (\bar{v}^2)^{1/2} = (3kT/2)^{1/2}$. If this function is known, all gas properties can be calculated by using it to obtain various averages. For example, the average momentum carried in a certain direction would give the viscosity. The velocity distribution for a gas at equilibrium was suggested by Maxwell in 1859 and is represented by the familiar bell-shaped curve that describes the normal, or Gaussian, distribution of random variables in large populations. Attempts to support more definitively this result and to extend it to nonequilibrium gases led to the formulation of the Boltzmann equation, which describes how collisions and external forces cause the velocity distribution to change. This equation is difficult to solve in any general sense, but some progress can be made by assuming that the deviations from the equilibrium distribution are small and are proportional to the external influences that cause the deviations, such as temperature, pressure, and composition differences. Even the resulting simpler equations remained unsolved for nearly 50 years until the work of Enskog and Chapman, with a single notable exception. The one case that was solvable dealt with molecules that interact with forces that fall off as the fifth power of their separation (*i.e.*, as $1/r^5$), for which Maxwell found an exact solution. Unfortunately, thermal diffusion happens to be exactly zero for molecules subject to this force law, so that phenomenon was missed.

It was later discovered that it is possible to use the solutions for the $1/r^5$ Maxwell model as a starting point and then calculate successive corrections for more general interactions. Although the calculations quickly increase in complexity, the improvement in accuracy is rapid, unlike the persistence-of-velocities corrections applied in mean free path theory. This refined version of kinetic theory is now highly developed, but it is quite mathematical and is not described here.

DEVIATIONS FROM THE IDEAL MODEL

Deviations from ideal gas behaviour occur both at low densities, where molecule-surface collisions become important, and at high densities, where a description in terms of only two-body collisions becomes inadequate. The low-density case can be handled in principle by including both molecule-surface and molecule-molecule collisions in the Boltzmann equation. Since this branch of the subject is now quite advanced and mathematical in character, only the high-density case will be discussed here.

Equation of state. To a first approximation, molecule-molecule collisions do not affect the ideal gas equation of state, $pV = RT$, but real gases at nonzero densities show deviations from this equation that are due to interactions among the molecules. Ever since the great advance made by van der Waals in 1873, an accurate universal formula relating p , v , and T has been sought. No completely satisfactory equation of state has been found, though important advances occurred in the 1970s and '80s. The only rigorous theoretical result available is an infinite-series expansion in powers of $1/v$, known as the virial equation of state:

$$\frac{pv}{RT} = 1 + \frac{B(T)}{v} + \frac{C(T)}{v^2} + \dots, \quad (28)$$

where $B(T)$, $C(T)$, \dots are called the second, third, \dots virial coefficients and depend only on the temperature and the particular gas. The virtue of this equation is that there is a rigorous connection between the virial coefficients and intermolecular forces, and experimental values of $B(T)$ were an early source (and still a useful one) of quantitative information on intermolecular forces. The drawback of the virial equation of state is that it is an infinite series

Velocity
distribution
function

Virial
equation of
state

and becomes essentially useless at high densities, which in practice are those greater than about the critical density. Also, the equation is wanting in that it does not predict condensation.

The most practical approaches to the equation of state for real fluids remain the versions of the principle of corresponding states first proposed by van der Waals.

Transport properties. Despite many attempts, there is still no satisfactory theory of the transport properties of dense fluids. Even the extension of the Boltzmann equation to include collisions of more than two bodies is not entirely clear. An important advance was made in 1921 by Enskog, but it is restricted to hard spheres and has not been extended to real molecules except in an empirical way to fit experimental measurements.

Attempts to develop a virial type of expansion in $1/v$ for the transport coefficients have failed in a surprising way. A formal theory was formulated, but, when the virial coefficients were evaluated for the tractable case of hard spheres, an infinite result was obtained for the coefficient of the $1/v^2$ term. This is a signal that a virial expansion is not accurate in a mathematical sense, and subsequent research showed that the error arose from a neglected term of the form $(1/v^2)\ln(1/v)$. It remains unknown how many similar problematic mathematical terms exist in the theory. Transport coefficients of dense fluids are usually described by some empirical extension of the Enskog hard-sphere theory or more commonly by some version of a principle of corresponding states. Much work clearly remains to be done. (E.A.M.)

PLASMA STATE

The plasma state of matter is unique like the solid, liquid, and gaseous states and is often considered the fourth state of matter. A plasma is an electrically conducting medium in which there are nearly equal numbers of positive and negative charges. The negative charge is usually carried by electrons, each of which has one unit of negative charge. The positive charge is typically carried by atoms or molecules that are missing those same electrons. In some rare but interesting cases, electrons missing from one type of atom or molecule become attached to another component, resulting in a plasma containing both positive and negative ions. The most extreme case of this type occurs when small but macroscopic dust particles become charged in a state referred to as a dusty plasma. The uniqueness of the plasma state is due to the importance of electric and magnetic forces that act on a plasma in addition to such forces as gravity that affect all forms of matter. Since these electromagnetic forces can act at large distances, a plasma will act collectively much like a fluid even when the particles seldom collide with one another.

Nearly all the visible matter in the universe exists in the plasma state, occurring predominantly in this form in the Sun and stars and in interplanetary and interstellar space. Auroras, lightning, and welding arcs are also plasmas; plasmas exist in neon and fluorescent tubes, in the crystal structure of metallic solids, and in many other phenomena and objects. The Earth itself is immersed in a tenuous plasma called the solar wind and is surrounded by a dense plasma called the ionosphere.

A plasma may be produced in the laboratory by heating a gas to an extremely high temperature, which causes such vigorous collisions between its atoms and molecules that electrons are ripped free, yielding the requisite electrons and ions. A similar process occurs inside stars. In space the dominant plasma formation process is photoionization, wherein photons from sunlight or starlight are absorbed by an existing gas, causing electrons to be emitted. Since the Sun and stars shine continuously, virtually all the matter becomes ionized in such cases, and the plasma is said to be fully ionized. This need not be the case, however, for a plasma may be only partially ionized. A completely ionized hydrogen plasma, consisting solely of electrons and protons (hydrogen nuclei), is the most elementary plasma.

The modern concept of the plasma state is of recent origin, dating back only to the early 1950s. Its history is interwoven with many disciplines. Three basic fields of study made unique early contributions to the development of plasma physics as a discipline: electric discharges, magnetohydrodynamics (in which a conducting fluid such as mercury is studied), and kinetic theory.

Interest in electric-discharge phenomena may be traced back to the beginning of the 18th century, with three English physicists—Michael Faraday in the 1830s and Joseph John Thomson and John Sealy Edward Townsend at the turn of the 19th century—laying the foundations of the present understanding of the phenomena. Irving Langmuir introduced the term plasma in 1923 while investigating electric discharges. In 1929 he and Lewi Tonks, another physicist working in the United States, used the term to

designate those regions of a discharge in which certain periodic variations of the negatively charged electrons could occur. They called these oscillations plasma oscillations, their behaviour suggesting that of a jellylike substance. Not until 1952, however, when two other American physicists, David Bohm and David Pines, first considered the collective behaviour of electrons in metals as distinct from that in ionized gases, was the general applicability of the concept of a plasma fully appreciated.

The collective behaviour of charged particles in magnetic fields and the concept of a conducting fluid are implicit in magnetohydrodynamic studies, the foundations of which were laid in the early and middle 1800s by Faraday and André-Marie Ampère of France. Not until the 1930s, however, when new solar and geophysical phenomena were being discovered, were many of the basic problems of the mutual interaction between ionized gases and magnetic fields considered. In 1942 Hannes Alfvén, a Swedish physicist, introduced the concept of magnetohydrodynamic waves. This contribution, along with his further studies of space plasmas, led to Alfvén's receipt of the Nobel Prize for Physics in 1970.

These two separate approaches—the study of electric discharges and the study of the behaviour of conducting fluids in magnetic fields—were unified by the introduction of the kinetic theory of the plasma state. This theory states that plasma, like gas, consists of particles in random motion, whose interactions can be through long-range electromagnetic forces as well as via collisions. In 1905 the Dutch physicist Hendrik Antoon Lorentz applied the kinetic equation for atoms (the formulation by the Austrian physicist Ludwig Eduard Boltzmann) to the behaviour of electrons in metals. Various physicists and mathematicians in the 1930s and '40s further developed the plasma kinetic theory to a high degree of sophistication. Since the early 1950s interest has increasingly focused on the plasma state itself. Space exploration, the development of electronic devices, a growing awareness of the importance of magnetic fields in astrophysical phenomena, and the quest for controlled thermonuclear (nuclear fusion) power reactors all have stimulated such interest. Many problems remain unsolved in space plasma physics research, owing to the complexity of the phenomena. For example, descriptions of the solar wind must include not only equations dealing with the effects of gravity, temperature, and pressure as needed in atmospheric science but also the equations of the Scottish physicist James Clerk Maxwell, which are needed to describe the electromagnetic field.

Just as a lightweight cork in water will bob up and down about its rest position, any general displacement of light electrons as a group with respect to the positive ions in a plasma leads to the oscillation of the electrons as a whole about an equilibrium state. In the case of the cork, the restoring force is provided by gravity; in plasma oscillations, it is provided by the electric force. These movements are the plasma oscillations that were studied by Langmuir and Tonks. Analogously, just as buoyancy effects guide water waves, plasma oscillations are related to waves in the electron component of the plasma called

Kinetic
theory

Langmuir waves. Wavelike phenomena play a critical role in the behaviour of plasmas.

The time τ required for an oscillation of this type is the most important temporal parameter in a plasma. The main spatial parameter is the Debye length, h , which is the distance traveled by the average thermal electron in time $\tau/2\pi$. A plasma can be defined in terms of these parameters as a partially or fully ionized gas that satisfies the following criteria: (1) a constituent electron may complete many plasma oscillations before it collides with either an ion or one of the other heavy constituents, (2) inside each sphere with a radius equal to the Debye length, there are many particles, and (3) the plasma itself is much larger than the Debye length in every dimension.

Another important temporal parameter is the time between collisions of particles. In any gas, separate collision frequencies are defined for collisions between all different particle types. The total collision frequency for a particular species is the weighted sum of all the separate frequencies.

Elastic and inelastic collisions

Two basic types of collision may occur: elastic and inelastic. In an elastic collision, the total kinetic energy of all the particles participating in the collision is the same before and after the event. In an inelastic collision, a fraction of the kinetic energy is transferred to the internal energy of the colliding particles. In an atom, for example, the electrons have certain allowed (discrete) energies and are said to be bound. During a collision, a bound electron may be excited—that is, raised from a low to a high energy state. This can occur, however, only by the expenditure of kinetic energy and only if the kinetic energy exceeds the difference between the two energy states. If the energy is sufficient, a bound electron may be excited to such a high level that it becomes a free electron, and the atom is said to be ionized; the minimum, or threshold, energy required to free an electron is called the ionization energy. Inelastic collisions may also occur with positive ions unless all the electrons have been stripped away. In general, only collisions of electrons and photons (quanta of electromagnetic radiation) with atoms and ions are significant in these inelastic collisions; ionization by a photon is called photoionization.

A molecule has additional discrete energy states, which may be excited by particle or photon collisions. At sufficiently high energies of interaction, the molecule can dissociate into atoms or into atoms and atomic ions. As in the case of atoms, collision of electrons and photons with molecules may cause ionization, producing molecular ions. In general, the reaction rate for inelastic collisions is similar to that of chemical reactions. At sufficiently high temperatures, the atoms are stripped of all electrons and become bare atomic nuclei. Finally, at temperatures of about 1,000,000 K or greater, nuclear reactions can occur—another form of inelastic collisions. When such reactions lead to the formation of heavier elements, the process is called thermonuclear fusion; mass is transmuted, and kinetic energy is gained instead of lost.

All sources of energy now existing on the Earth can be traced in one way or another to the nuclear fusion reactions inside the Sun or some long-extinct star. In such energy sources, gravity controls and confines the fusion process. The high temperatures required for the nuclear fusion reactions that take place in a hydrogen, or thermonuclear, bomb are attained by first igniting an atomic bomb, which produces a fission chain reaction. One of the great challenges of humankind is to create these high temperatures in a controlled manner and to harness the energy of nuclear fusion. This is the great practical goal of plasma physics—to produce nuclear fusion on the Earth. Confinement schemes devised by scientists use magnetic fields or the inertia of an implosion to guide and control the hot plasma.

Basic plasma physics

PLASMA FORMATION

Apart from solid-state plasmas, such as those in metallic crystals, plasmas do not usually occur naturally at the surface of the Earth. For laboratory experiments and technological applications, plasmas therefore must be pro-

duced artificially. Because the atoms of such alkalis as potassium, sodium, and cesium possess low ionization energies, plasmas may be produced from these by the direct application of heat at temperatures of about 3,000 K. In most gases, however, before any significant degree of ionization is achieved, temperatures in the neighbourhood of 10,000 K are required. A convenient unit for measuring temperature in the study of plasmas is the electron volt (eV), which is the energy gained by an electron in vacuum when it is accelerated across one volt of electric potential. The temperature, W , measured in electron volts is given by $W = T/12,000$ when T is expressed in kelvins. The temperatures required for self-ionization thus range from 2.5 to 8 electron volts, since such values are typical of the energy needed to remove one electron from an atom or molecule.

Because all substances melt at temperatures far below that level, no container yet built can withstand an external application of the heat necessary to form a plasma; therefore, any heating must be supplied internally. One technique is to apply an electric field to the gas to accelerate and scatter any free electrons, thereby heating the plasma. This type of ohmic heating is similar to the method in which free electrons in the heating element of an electric oven heat the coil. Because of their small energy loss in elastic collisions, electrons can be raised to much higher temperatures than other particles. For plasma formation a sufficiently high electric field must be applied, its exact value depending on geometry and the gas pressure. The electric field may be set up via electrodes or by transformer action, in which the electric field is induced by a changing magnetic field. Laboratory temperatures of about 10,000,000 K, or 8 kiloelectron volts (keV), with electron densities of about 10^{19} per cubic metre have been achieved by the transformer method. The temperature is eventually limited by energy losses to the outside environment. Extremely high temperatures, but relatively low-density plasmas, have been produced by the separate injection of ions and electrons into a mirror system (a plasma device using a particular arrangement of magnetic fields for containment). Other methods have used the high temperatures that develop behind a wave that is moving much faster than sound to produce what is called a shock front; lasers have also been employed.

Natural plasma heating and ionization occur in analogous ways. In a lightning-induced plasma, the electric current carried by the stroke heats the atmosphere in the same manner as in the ohmic heating technique described above. In solar and stellar plasmas the heating is internal and caused by nuclear fusion reactions. In the solar corona, the heating occurs because of waves that propagate from the surface into the Sun's atmosphere, heating the plasma much like shock-wave heating in laboratory plasmas. In the ionosphere, ionization is accomplished not through heating of the plasma but rather by the flux of energetic photons from the Sun. Far-ultraviolet rays and X rays from the Sun have enough energy to ionize atoms in the Earth's atmosphere. Some of the energy also goes into heating the gas, with the result that the upper atmosphere, called the thermosphere, is quite hot. These processes protect the Earth from energetic photons such as the ozone layer protects terrestrial life-forms from lower-energy ultraviolet light. The typical temperature 300 kilometres above the Earth's surface is 1,200 K, or about 0.1 eV. Although it is quite warm compared with the surface of the Earth, this temperature is too low to create self-ionization. When the Sun sets with respect to the ionosphere, the source of ionization ceases, and the lower portion of the ionosphere reverts to its nonplasma state. Some ions, in particular singly charged oxygen (O^+), live long enough that some plasma remains until the next sunrise. In the case of an aurora, a plasma is created in the nighttime or daytime atmosphere when beams of electrons are accelerated to hundreds or thousands of electron volts and smash into the atmosphere.

METHODS OF DESCRIBING PLASMA PHENOMENA

The behaviour of a plasma may be described at different levels. If collisions are relatively infrequent, it is useful

Artificial production of plasmas

to consider the motions of individual particles. In most plasmas of interest, a magnetic field exerts a force on a charged particle only if the particle is moving, the force being at right angles to both the direction of the field and the direction of particle motion. In a uniform magnetic field (B), a charged particle gyrates about a line of force. The centre of the orbit is called the guiding centre. The particle may also have a component of velocity parallel to the magnetic field and so traces out a helix in a uniform magnetic field. If a uniform electric field (E) is applied at right angles to the direction of the magnetic field, the guiding centre drifts with a uniform velocity of magnitude equal to the ratio of the electric to the magnetic field (E/B), at right angles to both the electric and magnetic fields. A particle starting from rest in such fields follows the same cycloidal path a dot on the rim of a rolling wheel follows. Although the "wheel" radius and its sense of rotation vary for different particles, the guiding centre moves at the same E/B velocity, independent of the particle's charge and mass. Should the electric field change with time, the problem would become even more complex. If, however, such an alternating electric field varies at the same frequency as the cyclotron frequency (*i.e.*, the rate of gyration), the guiding centre will remain stationary, and the particle will be forced to travel in an ever-expanding orbit. This phenomenon is called cyclotron resonance and is the basis of the cyclotron particle accelerator.

Cyclotron
resonance

The motion of a particle about its guiding centre constitutes a circular current. As such, the motion produces a dipole magnetic field not unlike that produced by a simple bar magnet. Thus, a moving charge not only interacts with magnetic fields but also produces them. The direction of the magnetic field produced by a moving particle, however, depends both on whether the particle is positively or negatively charged and on the direction of its motion. If the motion of the charged particles is completely random, the net associated magnetic field is zero. On the other hand, if charges of different sign have an average relative velocity (*i.e.*, if an electric current flows), then a net magnetic field over and above any externally applied field exists. The magnetic interaction between charged particles is therefore of a collective, rather than of an individual, particle nature.

At a higher level of description than that of the single particle, kinetic equations of the Boltzmann type are used. Such equations essentially describe the behaviour of those particles about a point in a small-volume element, the particle velocities lying within a small range about a given value. The interactions with all other velocity groups, volume elements, and any externally applied electric and magnetic fields are taken into account. In many cases, equations of a fluid type may be derived from the kinetic equations; they express the conservation of mass, momentum, and energy per unit volume, with one such set of equations for each particle type.

DETERMINATION OF PLASMA VARIABLES

The basic variables useful in the study of plasma are number densities, temperatures, electric and magnetic field strengths, and particle velocities. In the laboratory and in space, both electrostatic (charged) and magnetic types of sensory devices called probes help determine the magnitudes of such variables. With the electrostatic probe, ion densities, electron and ion temperatures, and electrostatic potential differences can be determined. Small search coils and other types of magnetic probes yield values for the magnetic field; and from Maxwell's electromagnetic equations the current and charge densities and the induced component of the electric field may be found. Interplanetary spacecraft have carried such probes to nearly every planet in the solar system, revealing to scientists such plasma phenomena as lightning on Jupiter and the sounds of Saturn's rings and radiation belts. In the early 1990s, signals were being relayed to the Earth from several spacecraft approaching the edge of the plasma boundary to the solar system, the heliopause.

In the laboratory the absorption, scattering, and excitation of neutral and high-energy ion beams are helpful in determining electron temperatures and densities; in gen-

eral, the refraction, reflection, absorption, scattering, and interference of electromagnetic waves also provide ways to determine these same variables. This technique has also been employed to remotely measure the properties of the plasmas in the near-space regions of the Earth using the incoherent scatter radar method. The largest single antenna is at the National Astronomy and Ionosphere Center at Arecibo in Puerto Rico. It has a circumference of 305 metres and was completed in 1963. It is still used to probe space plasmas to distances of 3,000 kilometres. The method works by bouncing radio waves from small irregularities in the electron gas that occur owing to random thermal motions of the particles. The returning signal is shifted slightly from the transmitted one—because of the Doppler-shift effect—and the velocity of the plasma can be determined in a manner similar to the way in which the police detect a speeding car. Using this method, the wind speed in space can be found, along with the temperature, density, electric field, and even the types of ions present. In geospace the appropriate radar frequencies are in the range of 50 to 1,000 megahertz (MHz), while in the laboratory, where the plasma densities and plasma frequencies are higher, microwaves and lasers must be used.

Aside from the above methods, much can be learned from the radiation generated and emitted by the plasma itself; in fact, this is the only means of studying cosmic plasma beyond the solar system. The various spectroscopic techniques covering the entire continuous radiation spectrum determine temperatures and identify such nonthermal sources as those pulses producing synchrotron radiations.

WAVES IN PLASMAS

The waves most familiar to people are the buoyancy waves that propagate on the surfaces of lakes and oceans and break onto the world's beaches. Equally familiar, although not necessarily recognized as waves, are the disturbances in the atmosphere that create what is referred to as the weather. Wave phenomena are particularly important in the behaviour of plasmas. In fact, one of the three criteria for the existence of a plasma is that the particle-particle collision rate be less than the plasma-oscillation frequency. This in turn implies that the collective interactions that control the plasma gas depend on the electric and magnetic field effects as much as, or more so than, simple collisions. Since waves are able to propagate, the possibility exists for force fields to act at large distances from the point where they originated.

Ordinary fluids can support the propagation of sound (acoustic) waves, which involve pressure, temperature, and velocity variations. Electromagnetic waves can propagate even in a vacuum but are slowed down in most cases by the interaction of the electric fields in the waves with the charged particles bound in the atoms or molecules of the gas. Although it is important for a complete description of electromagnetic waves, such an interaction is not very strong. In a plasma, however, the particles react in concert with any electromagnetic field (*e.g.*, as in an electromagnetic wave) as well as with any pressure or velocity field (*e.g.*, as in a sound wave). In fact, in a plasma sound wave the electrons and ions become slightly separated owing to their difference in mass, and an electric field builds up to bring them back together. The result is called an ion acoustic wave. This is just one of the many types of waves that can exist in a plasma. The brief discussion that follows touches on the main types in order of increasing wave-oscillation frequency.

At the lowest frequency are Alfvén waves, which require the presence of a magnetic field to exist. In fact, except for ion acoustic waves, the existence of a background magnetic field is required for any wave with a frequency less than the plasma frequency to occur in a plasma. Most natural plasmas are threaded by a magnetic field, and laboratory plasmas often use a magnetic field for confinement, so this requirement is usually met, and all types of waves can occur.

Alfvén
waves

Alfvén waves are analogous to the waves that occur on the stretched string of a guitar. In this case, the string represents a magnetic field line. When a small magnetic field disturbance takes place, the field is bent slightly, and

the disturbance propagates in the direction of the magnetic field. Since any changing magnetic field creates an electric field, an electromagnetic wave results. Such waves are the slowest and have the lowest frequencies of any known electromagnetic waves. For example, the solar wind streams out from the Sun with a speed greater than either electromagnetic (Alfvén) or sound waves. This means that, when the solar wind hits the Earth's outermost magnetic field lines, a shock wave results to "inform" the incoming plasma that an obstacle exists, much like the shock wave associated with a supersonic airplane. The shock wave travels toward the Sun at the same speed but in the opposite direction as the solar wind, so it appears to stand still with respect to the Earth. Because there are almost no particle-particle collisions, this type of collisionless shock wave is of great interest to space plasma physicists who postulate that similar shocks occur around supernovas and in other astrophysical plasmas. On the Earth's side of the shock wave, the heated and slowed solar wind interacts with the Earth's atmosphere via Alfvén waves propagating along the magnetic field lines.

The turbulent surface of the Sun radiates large-amplitude Alfvén waves, which are thought to be responsible for heating the corona to 1,000,000 K. Such waves can also produce fluctuations in the solar wind, and, as they propagate through it to the Earth, they seem to control the occurrence of magnetic storms and auroras that are capable of disrupting communication systems and power grids on the planet.

Types
of wave
motion

Two fundamental types of wave motion can occur: longitudinal, like a sound or ion acoustic wave, in which particle oscillation is in a direction parallel to the direction of wave propagation; and transverse, like a surface water wave, in which particle oscillation is in a plane perpendicular to the direction of wave propagation. In all cases, a wave may be characterized by a speed of propagation (u), a wavelength (λ), and a frequency (ν) related by an expression in which the velocity is equal to the product of the wavelength and frequency, namely, $u = \lambda\nu$. The Alfvén wave is a transverse wave and propagates with a velocity that depends on the particle density and the magnetic field strength. The velocity is equal to the magnetic flux density (B) divided by the square root of the mass density (ρ) times the permeability of free space (μ_0)—that is to say, $B/\sqrt{\mu_0\rho}$. The ion acoustic wave is a longitudinal wave and also propagates parallel to the magnetic field at a speed roughly equal to the average thermal velocity of the ions. Perpendicular to the magnetic field a different type of longitudinal wave called a magnetosonic wave can occur.

In these waves the plasma behaves as a whole, and the velocity is independent of wave frequency. At higher frequencies, however, the separate behaviour of ions and electrons causes the wave velocities to vary with direction and frequency. The Alfvén wave splits into two components, referred to as the fast and slow Alfvén waves, which propagate at different frequency-dependent speeds. At still higher frequencies these two waves (called the electron cyclotron and ion cyclotron waves, respectively) cause electron and cyclotron resonances (synchronization) at the appropriate resonance frequencies. Beyond these resonances, transverse wave propagation does not occur at all until frequencies comparable to and above the plasma frequency are reached.

At frequencies between the ion and electron gyrofrequencies lies a wave mode called a whistler. This name comes from the study of plasma waves generated by lightning. When early researchers listened to natural radio waves by attaching an antenna to an audio amplifier, they heard a strange whistling sound. The whistle occurs when the electrical signal from lightning in one hemisphere travels along the Earth's magnetic field lines to the other hemisphere. The trip is so long that some waves (those at higher frequencies) arrive first, resulting in the generation of a whistle-like sound. These natural waves were used to probe the region of space around the Earth before spacecraft became available. Such a frequency-dependent wave velocity is called wave dispersion because the various frequencies disperse with distance.

The speed of an ion acoustic wave also becomes dis-

persive at high frequencies, and a resonance similar to electron plasma oscillations occurs at a frequency determined by electrostatic oscillations of the ions. Beyond this frequency no sonic wave propagates parallel to a magnetic field until the frequency reaches the plasma frequency, above which electroacoustic waves occur. The wavelength of these waves at the critical frequency (ω_p) is infinite, the electron behaviour at this frequency taking the form of the plasma oscillations of Langmuir and Tonks. Even without particle collisions, waves shorter than the Debye length are heavily damped—*i.e.*, their amplitude decreases rapidly with time. This phenomenon, called Landau damping, arises because some electrons have the same velocity as the wave. As they move with the wave, they are accelerated much like a surfer on a water wave and thus extract energy from the wave, damping it in the process.

Magnetic fields are used to contain high-density, high-temperature plasmas because such fields exert pressures and tensile forces on the plasma. An equilibrium configuration is reached only when at all points in the plasma these pressures and tensions exactly balance the pressure from the motion of the particles. A well-known example of this is the pinch effect observed in specially designed equipment. If an external electric current is imposed on a cylindrically shaped plasma and flows parallel to the plasma axis, the magnetic forces act inward and cause the plasma to constrict, or pinch. An equilibrium condition is reached in which the temperature is proportional to the square of the electric current. This result suggests that any temperature may be achieved by making the electric current sufficiently large, the heating resulting from currents and compression. In practice, however, since no plasma can be infinitely long, serious energy losses occur at the ends of the cylinder; also, major instabilities develop in such a simple configuration. Suppression of such instabilities has been one of the major efforts in laboratory plasma physics and in the quest to control the nuclear fusion reaction.

A useful way of describing the confinement of a plasma by a magnetic field is by measuring containment time (τ_c), or the average time for a charged particle to diffuse out of the plasma; this time is different for each type of configuration. Various types of instabilities can occur in plasma. These lead to a loss of plasma and a catastrophic decrease in containment time. The most important of these is called magnetohydrodynamic instability. Although an equilibrium state may exist, it may not correspond to the lowest possible energy. The plasma, therefore, seeks a state of lower potential energy, just as a ball at rest on top of a hill (representing an equilibrium state) rolls down to the bottom if perturbed; the lower energy state of the plasma corresponds to a ball at the bottom of a valley. In seeking the lower energy state, turbulence develops, leading to enhanced diffusion, increased electrical resistivity, and large heat losses. In toroidal geometry, circular plasma currents must be kept below a critical value called the Kruskal-Shafranov limit, otherwise a particularly violent instability consisting of a series of kinks may occur. Although a completely stable system appears to be virtually impossible, considerable progress has been made in devising systems that eliminate the major instabilities. Temperatures on the order of 10,000,000 K at densities of 10^{19} particles per cubic metre and containment times as high as $1/50$ of a second have been achieved.

Containment

Kruskal-Shafranov limit

APPLICATIONS OF PLASMAS

The most important practical applications of plasmas lie in the future, largely in the field of power production. The major method of generating electric power has been to use heat sources to convert water to steam, which drives turbogenerators. Such heat sources depend on the combustion of fossil fuels, such as coal, oil, and natural gas, and fission processes in nuclear reactors. A potential source of heat might be supplied by a fusion reactor, with a basic element of deuterium-tritium plasma; nuclear fusion collisions between those isotopes of hydrogen would release large amounts of energy to the kinetic energy of the reaction products (the neutrons and the nuclei of hydrogen and helium atoms). By absorbing those products in a

surrounding medium, a powerful heat source could be created. To realize a net power output from such a generating station—allowing for plasma radiation and particle losses and for the somewhat inefficient conversion of heat to electricity—plasma temperatures of about 100,000,000 K and a product of particle density times containment time of about 10^{20} seconds per cubic metre are necessary. For example, at a density of 10^{20} particles per metre cubed, the containment time must be one second. Such figures are yet to be reached, although there has been much progress.

In general, there are two basic methods of eliminating or minimizing end losses from an artificially created plasma: the production of toroidal plasmas and the use of magnetic mirrors (see ATOMS: *Nuclear fusion*). A toroidal plasma is essentially one in which a plasma of cylindrical cross section is bent in a circle so as to close on itself. For such plasmas to be in equilibrium and stable, however, special magnetic fields are required, the largest component of which is a circular field parallel to the axis of the plasma. In addition, a number of turbulent plasma processes must be controlled to keep the system stable. In 1991 a machine called the JET (Joint European Torus) achieved the state known as ignition, a first step toward a practical fusion device.

Besides generating power, a fusion reactor might desalinate seawater. Approximately two-thirds of the world's land surface is uninhabited, with one-half of this area being arid. The use of both giant fission and fusion reactors in the large-scale evaporation of seawater could make irrigation of such areas economically feasible. Another possibility in power production is the elimination of the heat-steam-mechanical energy chain. One suggestion depends on the dynamo effect. If a plasma moves perpendicular to a magnetic field, an electromotive force, according to Faraday's law, is generated in a direction perpendicular to both the direction of flow of the plasma and the magnetic field. This dynamo effect can drive a current in an external circuit connected to electrodes in the plasma, and thus electric power may be produced without the need for steam-driven rotating machinery. This process is referred to as magnetohydrodynamic (MHD) power generation and has been proposed as a method of extracting power from certain types of fission reactors. Such a generator powers the auroras as the Earth's magnetic field lines tap electrical current from the MHD generator in the solar wind.

The inverse of the dynamo effect, called the motor effect, may be used to accelerate plasma. By pulsing cusp-shaped magnetic fields in a plasma, for example, it is possible to achieve thrusts proportional to the square of the magnetic field. Motors based on such a technique have been proposed for the propulsion of craft in deep space. They have the advantage of being capable of achieving large exhaust velocities, thus minimizing the amount of fuel carried.

A practical application of plasma involves the glow discharge that occurs between two electrodes at pressures of one-thousandth of an atmosphere or thereabouts. Such glow discharges are responsible for the light given off by neon tubes and such other light sources as fluorescent lamps, which operate by virtue of the plasmas they produce in electric discharge. The degree of ionization in such plasmas is usually low, but electron densities of 10^{16} to 10^{18} electrons per cubic metre can be achieved with an electron temperature of 100,000 K. The electrons responsible for current flow are produced by ionization in a region near the cathode, with most of the potential difference between the two electrodes occurring there. This region does not contain a plasma, but the region between it and the anode (*i.e.*, the positive electrode) does.

Other applications of the glow discharge include electronic switching devices; it and similar plasmas produced by radio-frequency techniques can be used to provide ions for particle accelerators and act as generators of laser beams. As the current is increased through a glow discharge, a stage is reached when the energy generated at the cathode is sufficient to provide all the conduction electrons directly from the cathode surface, rather than from gas between the electrodes. Under this condition the large cathode potential difference disappears, and the plasma column contracts. This new state of electric discharge is

called an arc. Compared with the glow discharge, it is a high-density plasma and will operate over a large range of pressures. Arcs are used as light sources for welding, in electronic switching, for rectification of alternating currents, and in high-temperature chemistry. Running an arc between concentric electrodes and injecting gas into such a region causes a hot, high-density plasma mixture called a plasma jet to be ejected. It has many chemical and metallurgical applications.

Natural plasmas

EXTRATERRESTRIAL FORMS

It has been suggested that the universe originated as a violent explosion about 10 billion years ago and initially consisted of a fireball of completely ionized hydrogen plasma. Irrespective of the truth of this, there is little matter in the universe now that does not exist in the plasma state. The observed stars are composed of plasmas, as are interstellar and interplanetary media and the outer atmospheres of planets. Scientific knowledge of the universe has come primarily from studies of electromagnetic radiation emitted by plasmas and transmitted through them and, since the 1960s, from space probes within the solar system.

In a star the plasma is bound together by gravitational forces, and the enormous energy it emits originates in thermonuclear fusion reactions within the interior. Heat is transferred from the interior to the exterior by radiation in the outer layers, where convection is of greater importance. In the vicinity of a hot star, the interstellar medium consists almost entirely of completely ionized hydrogen, ionized by the star's ultraviolet radiation. Such regions are referred to as H II regions. The greater proportion by far of interstellar medium, however, exists in the form of neutral hydrogen clouds referred to as H I regions. Because the heavy atoms in such clouds are ionized by ultraviolet radiation (or photoionized), they also are considered to be plasmas, although the degree of ionization is probably only one part in 10,000. Other components of the interstellar medium are grains of dust and cosmic rays, the latter consisting of very high-energy atomic nuclei completely stripped of electrons. The almost isotropic velocity distribution of the cosmic rays may stem from interactions with waves of the background plasma.

Throughout this universe of plasma there are magnetic fields. In interstellar space magnetic fields are about 5×10^{-6} gauss (a unit of magnetic field strength) and in interplanetary space 5×10^{-5} gauss, whereas in intergalactic space they could be as low as 10^{-9} gauss. These values are exceedingly small compared with the Earth's surface field of about 5×10^{-1} gauss. Although small in an absolute sense, these fields are nevertheless gigantic, considering the scales involved. For example, to simulate interstellar phenomena in the laboratory, fields of about 10^{15} gauss would be necessary. Thus, these fields play a major role in nearly all astrophysical phenomena. On the Sun the average surface field is in the vicinity of 1 to 2 gauss, but magnetic disturbances arise, such as sunspots, in which fields of between 10 and 1,000 gauss occur. Many other stars are also known to have magnetic fields.

Magnetic fields in space

Table 8: Various Natural Plasmas and Their Electron Densities and Temperatures

plasma	n_e (per cu m)	T_e (K)
Sun		
Centre	10^{31}	1.5×10^7
Photosphere	10^{20}	4,200
Chromosphere	10^{17} – 10^{20}	5×10^5
Corona	10^{13}	1.5×10^6
Solar wind (near Earth)	5×10^6	4×10^5
Interstellar space		
H II regions	10^6	10^4
H I regions	10^2	100–125
Intergalactic space	1	3?
Earth		
Outer magnetosphere	10^6 – 10^7	10^4
Plasmasphere	10^9 – 10^{10}	10^4
Ionosphere	10^{11} – 10^{12}	250–3,000
Metals	10^{28}	10^4

Current from the dynamo effect

Use in radio switching devices

Field strengths of 10^{-3} gauss are associated with various extragalactic nebulae from which synchrotron radiation has been observed.

SOLAR-TERRESTRIAL FORMS

The visible region of the Sun is the photosphere (see Table 8), with its radiation being about the same as the continuum radiation from a 5,000 K blackbody. Lying above the photosphere is the chromosphere, which is observed by the emission of line radiation from various atoms and ions. Outside the chromosphere, the corona expands into the ever-blowing solar wind (see below), which on passing through the planetary system eventually encounters the interstellar medium. The corona can be seen in spectacular fashion when the Moon eclipses the bright photosphere (see Figure 43). During the times in which sunspots are greatest in number (called the sunspot maximum), the corona is very extended and the solar wind is fierce. Sunspot activity waxes and wanes with roughly an 11-year cycle. During the mid-1600s and early 1700s, sunspots virtually disappeared for a period known as the Maunder minimum. This time coincided with the Little Ice Age in Europe, and much conjecture has arisen about the possible effect of sunspots on climate. Periodic variations similar to that of sunspots have been observed in tree rings and lake-bed sedimentation. If real, such an effect is important because it implies that the Earth's climate is fragile.

In 1958 the American astrophysicist Eugene Parker showed that the equations describing the flow of plasma in the Sun's gravitational field had one solution that allowed the gas to become supersonic and to escape the Sun's pull. The solution was much like the description of a rocket nozzle in which the constriction in the flow is analogous to the effect of gravity. Parker predicted the Sun's atmosphere would behave just as this particular solar-wind solution predicts rather than according to the solar-breeze solutions suggested by others. The interplanetary satellite probes of the 1960s proved his solution to be correct.

The solar wind is a collisionless plasma made up primarily of electrons and protons and carries an outflow of matter moving at supersonic and super-Alfvénic speed. The wind takes with it an extension of the Sun's magnetic

field, which is frozen into the highly conducting fluid. In the region of the Earth, the wind has an average speed of 400 kilometres per second; and, when it encounters the planet's magnetic field, a shock front develops, the pressures acting to compress the field on the side toward the Sun and elongate it on the nightside (in the Earth's lee away from the Sun). The Earth's magnetic field is therefore confined to a cavity called the magnetosphere, into which the direct entry of the solar wind is prohibited. This cavity extends for about 10 Earth radii on the Sun's side and about 1,000 Earth radii on the nightside.

Inside this vast magnetic field a region of circulating plasma is driven by the transfer of momentum from the solar wind. Plasma flows parallel to the solar wind on the edges of this region and back toward the Earth in its interior. The resulting system acts as a secondary magnetohydrodynamic generator (the primary one being the solar wind itself). Both generators produce potential on the order of 100,000 volts. The solar-wind potential appears across the polar caps of the Earth, while the magnetospheric potential appears across the auroral oval. The latter is the region of the Earth where energetic electrons and ions precipitate into the planet's atmosphere, creating a spectacular light show. This particle flux is energetic enough to act as a new source of plasmas even when the Sun is no longer shining. The auroral oval becomes a good conductor; and large electric currents flow along it, driven by the potential difference across the system. These currents commonly are on the order of 1,000,000 amperes.

The plasma inside the magnetosphere is extremely hot (1–10 million K) and very tenuous (1–10 particles per cubic centimetre). The particles are heated by a number of interesting plasma effects, the most curious of which is the auroral acceleration process itself. A particle accelerator that may be the prototype for cosmic accelerators throughout the universe is located roughly one Earth radius above the auroral oval and linked to it by all-important magnetic field lines. In this region the auroral electrons are boosted by a potential difference on the order of three to six kilovolts, most likely created by an electric field parallel to the magnetic field lines and directed away from the Earth. Such a field is difficult to explain because magnetic

Solar
wind

Stephen J. Edberg, cover photograph *Reviews of Geophysics*,
vol. 30, no. 1, published 1992 by the American Geophysical Union



Figure 43: The solar corona as seen from La Paz, Mex., during the July 11, 1991, total eclipse. The corona is the source of the solar wind plasma that continuously bathes the Earth.

field lines usually act like nearly perfect conductors. The auroras occur on magnetic field lines that—if it were not for the distortion of the Earth's dipole field—would cross the equatorial plane at a distance of 6–10 Earth radii.

Closer to the Earth, within about 4 Earth radii, the planet wrests control of the system away from the solar wind. Inside this region the plasma rotates with the Earth, just as its atmosphere rotates with it. This system can also be thought of as a magnetohydrodynamic generator in which the rotation of the atmosphere and the ionospheric plasma in it create an electric field that puts the inner magnetosphere in rotation about the Earth's axis. Since this inner region is in contact with the dayside of the Earth where the Sun creates copious amounts of plasma in the ionosphere, the inner zone fills up with dense, cool plasma to form the plasmasphere. On a planet such as Jupiter, which has both a larger magnetic field and a higher rotation rate than the Earth, planetary control extends much farther from the surface.

The ionosphere

At altitudes below about 2,000 kilometres, the plasma is referred to as the ionosphere. Thousands of rocket probes have helped chart the vertical structure of this region of the atmosphere, and numerous satellites have provided latitudinal and longitudinal information. The ionosphere was discovered in the early 1900s when radio waves were found to propagate "over the horizon." If radio waves have frequencies near or below the plasma frequency, they cannot propagate throughout the plasma of the ionosphere and thus do not escape into space; they are instead either reflected or absorbed. At night the absorption is low since little plasma exists at the height of roughly 100 kilometres where absorption is greatest. Thus, the ionosphere acts as an effective mirror, as does the Earth's surface, and waves can be reflected around the entire planet much as in a waveguide. A great communications revolution was initiated by the wireless, which relied on radio waves to transmit audio signals. Development continues to this day with satellite systems that must propagate through the ionospheric plasma. In this case, the wave frequency must be higher than the highest plasma frequency in the ionosphere so that the waves will not be reflected away from the Earth.

The dominant ion in the upper atmosphere is atomic oxygen, while below about 200 kilometres molecular oxygen and nitric oxide are most prevalent. Meteor showers also provide large numbers of metallic atoms of elements such as iron, silicon, and magnesium, which become ionized in sunlight and last for long periods of time. These form vast ion clouds, which are responsible for much of the fading in and out of radio stations at night.

Noctilucent clouds

A more normal type of cloud forms at the base of the Earth's plasma blanket in the summer polar mesosphere regions. Located at an altitude of 85 kilometres, such a cloud is the highest on Earth and can be seen only when darkness has just set in on the planet. Hence, clouds of this kind have been called noctilucent clouds. They are thought to be composed of charged and possibly dusty ice crystals that form in the coldest portion of the atmosphere at a temperature of 120 K. This unusual medium has much in common with dusty plasmas in planetary rings and other cosmic systems. Noctilucent clouds have been increasing in frequency throughout the 20th century and may be a forerunner of global change.

High-energy particles also exist in the magnetosphere. At about 1.5 and 3.5 Earth radii from the centre of the planet, two regions contain high-energy particles. These regions are the Van Allen radiation belts, named after the American scientist James Van Allen, who discovered them using radiation detectors aboard early spacecraft. The charged particles in the belts are trapped in the mirror system formed by the Earth's magnetic dipole field.

Plasma can exist briefly in the lowest regions of the Earth's atmosphere. In a lightning stroke an oxygen-nitrogen plasma is heated at approximately 20,000 K with an ionization of about 20 percent, similar to that of a laboratory arc. Although the stroke is only a few centimetres thick and lasts only a fraction of a second, tremendous energies are dissipated. A lightning flash between the ground and a cloud, on the average, consists of four such strokes

in rapid succession. At all times, lightning is occurring somewhere on the Earth, charging the surface negatively with respect to the ionosphere by roughly 200,000 volts, even far from the nearest thunderstorm. If lightning ceased everywhere for even one hour, the Earth would discharge. An associated phenomenon is ball lightning. There are authenticated reports of glowing, floating, stable balls of light several tens of centimetres in diameter occurring at times of intense electrical activity in the atmosphere. On contact with an object, these balls release large amounts of energy. Although lightning balls are probably plasmas, so far no adequate explanation of them has been given.

Considering the origins of plasma physics and the fact that the universe is little more than a vast sea of plasma, it is ironic that the only naturally occurring plasmas at the surface of the Earth besides lightning are those to be found in ordinary matter. The free electrons responsible for electrical conduction in a metal constitute a plasma. Ions are fixed in position at lattice points, and so plasma behaviour in metals is limited to such phenomena as plasma oscillations and electron cyclotron waves (called helicon waves) in which the electron component behaves separately from the ion component. In semiconductors, on the other hand, the current carriers are electrons and positive holes, the latter behaving in the material as free positive charges of finite mass. By proper preparation, the number of electrons and holes can be made approximately equal so that the full range of plasma behaviour can be observed.

The importance of magnetic fields in astrophysical phenomena has already been noted. It is believed that these

Lightning flashes and ball lightning

Frederick J. Rich, cover photograph. *Reviews of Geophysics*, vol. 28, no. 3, published 1990 by the American Geophysical Union



Figure 44: An image in visible light of the eastern United States and Canada, taken by the optical line scanner aboard the F9 satellite of the Defense Meteorological Satellite Program. The image was obtained on March 14, 1989, during a major magnetic storm. The white band stretching from Hudson Bay to Ohio is due to electron impact on the Earth's atmosphere that results in the creation of a dense plasma and the light emission called an aurora.

fields are produced by self-generating dynamos, although the exact details are still not fully understood. In the case of the Earth, differential rotation in its liquid conducting core causes the external magnetic dipole field (manifest as the North and South poles). Cyclonic turbulence in the liquid, generated by heat conduction and Coriolis forces (apparent forces accompanying all rotating systems, including the heavenly bodies), generates the dipole field from these loops. Over geologic time, the Earth's field occasionally becomes small and then changes direction, the North Pole becoming the South Pole and vice versa. During the times in which the magnetic field is small, cosmic rays can more easily reach the Earth's surface and may affect life forms by increasing the rate at which genetic mutations occur.

Similar magnetic-field generation processes are believed to occur in both the Sun and the Milky Way Galaxy.

In the Sun the circular internal magnetic field is made observable by lines of force apparently breaking the solar surface to form exposed loops; entry and departure points are what are observed as sunspots. Although the exterior magnetic field of the Earth is that of a dipole, this is further modified by currents in both the ionosphere and magnetosphere. Lunar and solar tides in the ionosphere lead to motions across the Earth's field that produce currents, like a dynamo, that modify the initial field. The auroral oval current systems discussed earlier create even larger magnetic-field fluctuations. The intensity of these currents is modulated by the intensity of the solar wind, which also induces or produces other currents in the magnetosphere. Such currents taken together constitute the essence of a magnetic storm. An image of the eastern United States and Canada taken during a major magnetic storm is shown in Figure 44. (B.S.L./M.C.K.)

CLUSTERS

Atoms and molecules are the smallest forms of matter typically encountered under normal conditions and are in that sense the basic building blocks of the material world. As described in the previous section, there are phenomena, such as lightning and electric discharges of other kinds, that allow free electrons to be observed, but these are exceptional occurrences. It is of course in its gaseous state that matter is encountered at its atomic or molecular level; in gases each molecule is an independent entity, only occasionally and briefly colliding with another molecule or with a confining wall.

In contrast to the free-molecule character of gases, the condensed phases of matter—as liquids, crystalline solids, and glasses are called—depend for their properties on the constant proximity of all their constituent atoms. The extent to which the identities of the molecular constituents are maintained varies widely in these condensed forms of matter. Weakly bound solids, such as solid carbon dioxide (dry ice), or their liquid counterparts, such as liquid air or liquid nitrogen, are made up of molecules whose properties differ only slightly from the properties of the same molecules in gaseous form: such solids or liquids are simply molecules packed tightly enough to be in constant contact. These are called van der Waals solids or liquids, after Johannes D. van der Waals, the Dutch physicist who described the weak forces that just manage to hold these materials together if they are cold enough. In other solids, like diamond, graphite, silicon, or quartz, the individual atoms retain their identity, but there are no identifiable molecules in their structures. The forces between the constituent atoms are roughly as strong as the forces that hold atoms together in the strongly bound covalent molecules that make up most common substances. Negatively charged electrons act as a "glue" to hold the positively charged nuclei together and are more or less confined to the vicinity of the so-called home-base nuclei with which they are associated: they are not free to roam through the entire solid or liquid. These materials are said to be covalently bound and are electrical insulators. They are best described as neutral atoms held together by covalent bonds and are essentially one giant molecule.

Another kind of bonding found in condensed matter is exhibited by sodium chloride, ordinary table salt, which is composed of positive sodium ions (Na^+) and negative chloride ions (Cl^-). Such ionic compounds are held together by the mutual attraction of the oppositely charged ions: because of their locations, these attractions are stronger than the repulsions of the ions with like charges. Each ion in an ionic crystal is surrounded by nearest neighbours of opposite charge. The consequence is that the binding energies of ionic compounds are large, comparable to those of strongly bound covalent substances.

Metallic bonding is another type of binding found in condensed matter. Electrons moving between the positive atomic cores (*i.e.*, the nuclei plus inner-shell, tightly bound electrons) form an electron cloud: the attractions between the positive cores and the negative charges that make up

the cloud hold metals together. Metals differ from covalently bound insulators in that those electrons responsible for the cohesion of the metals move freely throughout the metal when given the slightest extra energy. For example, under the influence of the electric field produced in a copper wire when its ends are connected to the terminals of a battery, electrons move through the wire from the end connected to the battery's negative pole toward the end connected to its positive pole. An electric field applied to a metal generates an electric current, but the same electric field applied to a covalent insulator does not (see below *Comparison with other forms of matter*). The net binding forces between electrons and atomic cores of a metal are comparable in strength to those that hold ionic compounds together.

As mentioned above, liquids constitute a condensed or dense phase of matter, but their atomic arrangement differs from that of solids. In a liquid the constituent atoms are only slightly farther apart than they are in a solid, but that small difference is significant enough to allow the atoms or molecules that constitute a liquid to move around and to assume a full range of geometric configurations. Atoms of the same kind can trade places and can wander through the liquid by the random-walk process called diffusion. In general, materials that can form solids can also form liquids, but some, such as carbon dioxide, can only enter the liquid state under excess pressure. At least one substance, helium, can form a liquid while having no known solid form.

Materials that form solids and liquids can exhibit another form, one that may be solidlike or liquidlike but that has properties somewhat different from those of the bulk. This is the form of matter consisting of exceedingly small particles that are called clusters. Clusters are aggregates of atoms, molecules, or ions that adhere together under forces like those that bind the atoms, ions, or molecules of bulk matter: because of the manner in which they are prepared, clusters remain as tiny particles at least during the course of an experiment. There are clusters held together by van der Waals forces, by ionic forces, by covalent bonds, and by metallic bonds. Despite the similarity of the forces that bind both clusters and bulk matter, one of the fascinating aspects of clusters is that their properties differ from those of the corresponding bulk material: that characteristic affords the opportunity to learn about the properties of bulk matter by studying how, as the number of constituent particles increases, the properties of clusters evolve into those of bulk matter. For example, a cluster of 20 or 30 atoms typically has a melting point far lower than that of the corresponding bulk. The electrical properties of clusters also differ in some instances from those of the bulk matter: clusters of only a few atoms of mercury are insulators, held together by weak van der Waals forces, but clusters of hundreds of mercury atoms are metallic. One of the puzzles posed by clusters is the question of how properties of small clusters evolve with size into properties of bulk matter.

Binding in condensed phases of matter

Diffusion of liquids

Comparison with other forms of matter

Clusters versus bulk matter

Several characteristics differentiate clusters from molecules and bulk matter. They differ from bulk matter, first and foremost, in size; whether three particles bound together constitute a cluster is a matter of choice and convention, but an aggregate of four or more atoms or molecules certainly comprises a cluster. Such a small cluster would differ markedly from bulk matter in almost all its properties. A second difference between clusters and bulk matter is the variability of the properties of clusters with the number of their constituent particles. The properties of a lump of bulk matter remain unchanged by the addition or subtraction of a few atoms or molecules, whereas the properties of a small cluster vary significantly and, in general, neither uniformly nor even in the same direction with a change in the number of constituent particles. Medium-size clusters have properties that vary smoothly with the number of constituent particles (denoted N), but their properties, such as the melting point, differ significantly from those of the corresponding bulk. Large clusters have properties that vary smoothly with N and clearly merge into those of their bulk counterparts. This distinction, while not extremely precise, is quite useful. For example, the average binding energies—that is, the average energy per constituent atom or molecule required to separate the particles from each other—vary widely with N for small clusters. The reason for this wide range is that clusters of certain values of N , known as magic numbers, can take on unusually stable geometric structures that yield large binding energies, while others with different small values of N have no especially stable forms and therefore only relatively low binding energies. The binding energies of medium-size clusters vary rather smoothly with N , but they are in general considerably lower than the binding energies of bulk matter. The most important reason for this trend is that in a body of bulk matter almost all the particles are in the interior, while in a cluster most of the particles are on the surface. In a cluster of 13 atoms of copper or argon, for example, 12 of the atoms are on the surface. In a cluster of 55 argon atoms, 42 atoms are on the surface, and, in a cluster of 137 argon atoms, 82 are on the surface. Surface atoms are bonded only to atoms in their own layer and to those directly beneath them, so they have fewer atoms holding them to the main body of matter, whether cluster or bulk, than do atoms in the interior. Hence, the average binding energies of atoms in clusters are normally considerably less than those of bulk matter.

An important difference between clusters, in particular small and medium-size clusters, and bulk solids is the structure that is assumed by their most stable form. Most bulk solids are crystalline. This means that their atomic structures consist of periodic lattices—*i.e.*, structures that repeat over and over so that every unit composed of a few neighbouring atoms is indistinguishable from other groups of atoms that have exactly the same arrangement. In a simple cubic crystal, for example, all the atoms lie at the corners of cubes (in fact at a point common to eight equivalent cubes), and all these lattice points are identical. Such structures are called periodic. Most clusters, by contrast, have structures that are not periodic; many have the form of icosahedrons, incomplete icosahedrons, or other polyhedral structures that cannot grow into periodic lattices. One of the challenging puzzles of cluster science is to explain how, as an aggregate grows, it transforms from a polyhedral cluster-type structure into a crystalline lattice-type structure.

Furthermore, some properties of clusters reflect their small size in more subtle ways that depend on quantum mechanical phenomena. These are generally much more pronounced in exceedingly small systems than in bulk or macroscopic samples. One such property is the nature of the energy levels occupied by the electrons. In a macroscopic sample the energies of the states available to an electron are, in principle, discrete but are merged into bands consisting of many energy levels. Within each band the intervals of energy between those discrete levels are too tiny to be discerned; only the gaps between the bands

are large enough to be important because they correspond to ranges of energy that are forbidden to the electrons. In fact, it is the contrast in the mobility of electrons that differentiates insulators from electrical conductors. In even a very cold metal, only an infinitesimal amount of excess energy is required to promote a few electrons into the previously empty energy levels in which they can move freely throughout the material. If an electric field is applied to the metal, the negatively charged electrons move toward the positive pole of the field so that a net current flows in the metal. It is the motion of these electrons, driven by an applied field, that makes metals conductors of electricity. In an insulator the electrons fill all the energy levels up to the top of the highest-energy occupied band. This means that at least the full energy of the forbidden interval, called the band gap, must be imparted to any electron to promote it to an allowed state where it may travel readily through the material. In an insulator this is far more energy than is normally available, and so no electrons are in states that allow them to move freely; such materials cannot conduct electric currents.

Clusters containing only a small number of metal atoms have so few available quantum states for their electrons that these states must be considered discrete, not as components of a dense band of available states. In this sense, small clusters of metal atoms are like conventional molecules rather than like bulk metals. Medium-size clusters of metal atoms have electronic energy states that are

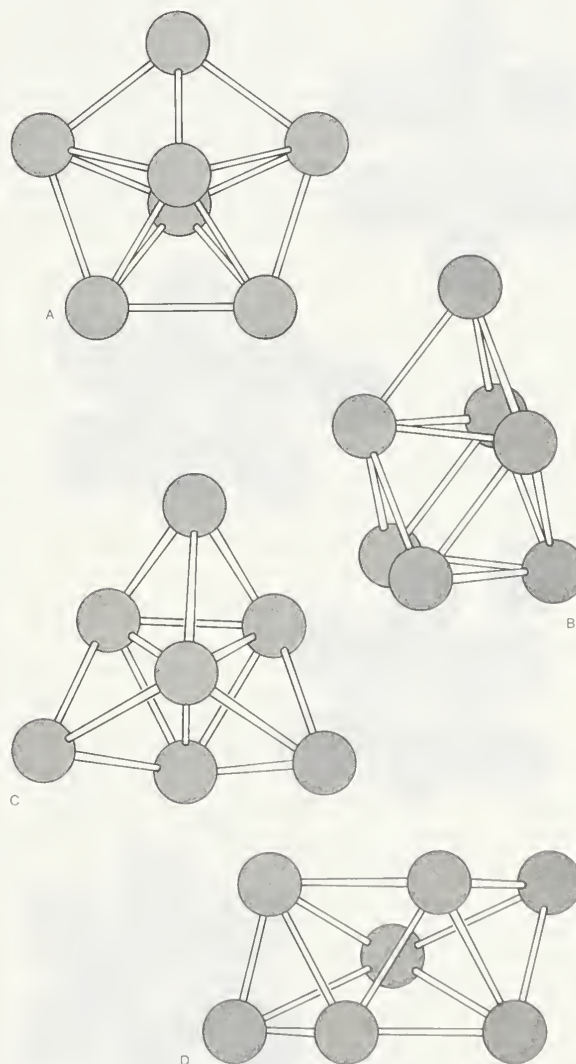


Figure 45: The four stable geometric structures of the seven-atom cluster of argon, in order of increasing energy: (A) A pentagonal bipyramid. (B) A regular octahedron with one face capped by the seventh atom. (C) A regular tetrahedron with three of its faces capped by other atoms. (D) A trigonal bipyramid with two of its faces capped by other atoms; although this has the highest energy of the four structures, it is very close in energy to the tricapped tetrahedron.

Quantum mechanical considerations

close enough together to be treated like the bands of bulk metals, but the conducting properties of these clusters are different from those of the bulk. Electrons driven by a constant electric field in a bulk metal can travel distances that are extremely long compared with atomic dimensions before they encounter any boundaries at the edges of the metal. Electrons in metal clusters encounter the boundaries of their cluster in a much shorter distance. Hence, metal clusters do not conduct electricity like bulk metals; if they are subjected to rapidly oscillating electric fields, such as those of visible, infrared, or microwave radiation, their "free" electrons are driven first one way and then back in the opposite direction over distances smaller than the dimensions of the cluster (see below *Physical properties*). If they are subjected to constant or low-frequency electric fields, such as the common 60-hertz fields that drive ordinary household currents, the electrons reach the boundaries of their clusters and can go no farther. Thus, the equivalent of conduction is not seen at low frequencies.

Clusters
versus
molecules

The manner in which clusters differ from molecules is more of a categorical nature than one of physical properties. Molecules have a definite composition and geometry; with few exceptions clusters can be made of any number of particles and may have any of several geometries. The four possible structures of a cluster of seven argon atoms are shown in Figure 45, and the lowest and next three higher-energy structures of a 13-atom cluster of argon are illustrated in Figure 46. The 13-atom cluster has the form

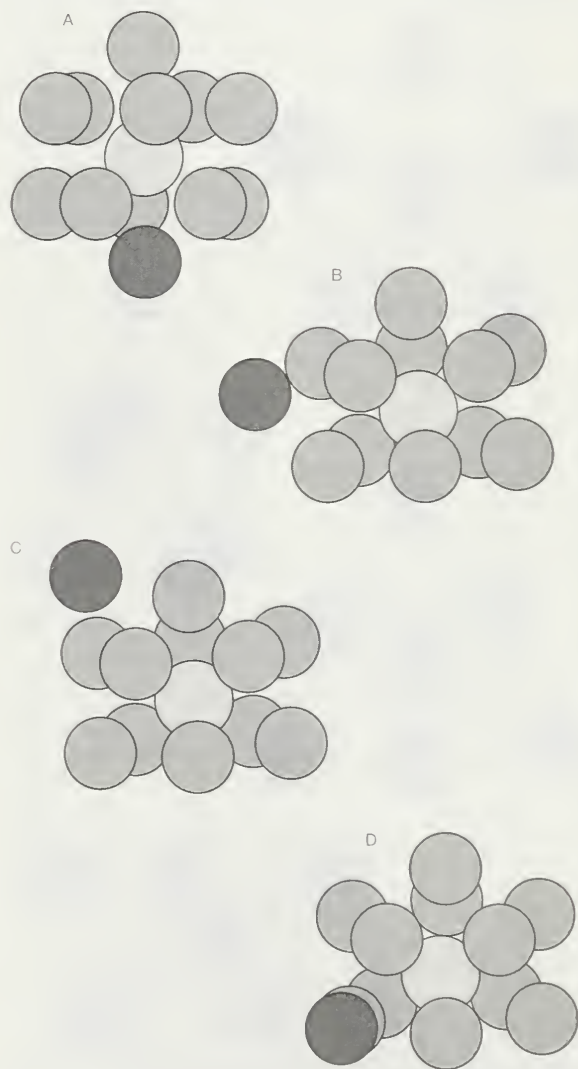


Figure 46: The four lowest-energy structures of the 13-atom cluster of argon: (A) The structure of lowest energy—the regular icosahedron of 12 atoms around a central atom. (B,C,D) The three structures, which have almost equal energy, formed by removing one of the 12 equivalent atoms from the 13-atom cluster in (A) from its shell and placing it into one of the three types of triangular faces in the resulting cluster.

of a regular icosahedron of 12 argon atoms around a central atom and is particularly stable.

Despite their multiplicity of structures, small clusters of fixed size, undergoing vibrations of small amplitude around a single geometry, are in most respects indistinguishable from molecules. If such clusters are given energy that is not great enough in magnitude to break them into separate parts, they may assume other geometries, alternating among these structural forms. This phenomenon is rarely seen with conventional molecules, but it is not unknown for energized molecules to exhibit more than one structure and to pass among them.

All in all, small clusters are much like molecules and are often considered to be molecules, while very large clusters are quite similar to bulk matter. The properties of clusters whose size is between these extremes may be like either or like neither.

Methods of study

Clusters can be studied by experiment, by theoretical analysis, and by simulation with computer-generated models. For several reasons they cannot be studied in the same manner as bulk matter. First, if individual clusters are allowed to coalesce into a mass, they will actually turn into bulk matter, so they must be kept separated. Second, it is desirable (but not always possible) to conduct experiments that distinguish the size and structure of each kind of cluster under observation. Because of these two considerations, experiments with clusters are usually more difficult than those with either specific molecules or bulk matter. Most of the difficulties arise from the same properties that make clusters interesting: the ease with which their sizes and compositions are varied and the variety of structures available for clusters of almost any given size.

Because of these difficulties, most experiments on clusters have been carried out with the clusters isolated in the gas phase; a few studies have been done with them in solution or in frozen matrices. Clusters can be prepared in the gas phase and then either studied in that form or captured into solvents or matrices or onto surfaces. They may be made by condensation of atoms or molecules or by direct blasting of matter from solids. In the most generally used method, a gas containing the gaseous cluster material is cooled by passing it under high pressure through a fine hole or slot. The expansion cools the gas rapidly from its initial temperature—usually room temperature but much higher if the cluster material is solid at room temperature—to a temperature not far above absolute zero. If, for example, argon gas is expanded in this way, it condenses into clusters if the pressure is not too high and the aperture is not too small; if the conditions are too extreme, the argon instead turns to snow and condenses.

Inert gases are often used as the medium by which other materials, in a gaseous or vaporous state, are transported from the ovens or other sources where they have been gasified and through the jets that cool them and turn them into clusters. One especially popular and interesting method in which solids are vaporized is by the action of intense laser beams on solid surfaces. Often called ablation, this process is an effective means of vaporizing even highly refractory materials like solid carbon. The ablated material is then carried through the cooling jet by an inert gas such as helium or argon.

Once the clusters have been formed, they can be studied in a variety of ways. One of the first techniques was simply to ionize the clusters, either with ultraviolet radiation (usually from a laser) or by electron impact. The gaseous ionized clusters are accelerated by an electric field and then analyzed according to their masses (see ANALYSIS AND MEASUREMENT: *Mass spectrometry*); these results immediately reveal the number of atoms or molecules in the cluster. The analysis yields the distribution of the relative abundances of clusters of different sizes in the beam. If the experiment is done with considerable care, the abundance distribution corresponds to the true relative stabilities of the clusters of different sizes. However, like many experiments with clusters, these can either provide results consistent with the equilibrium conditions that re-

Preparation
of clusters

flect those relative stabilities, or they can give results that reflect the rates of the cluster-forming processes rather than the equilibrium characteristics, as the latter may take far longer to reach than the time required to form clusters. Some of the implications of the abundances found in such experiments are discussed below in the section *Structure and properties: Structure*.

Because of the conditions under which clusters are formed, their distributions contain many different sizes and, in some instances, different shapes. Because chemists seek to characterize clusters of a single size and geometry, the clusters must first be sorted on that basis. If the clusters carry charge, they can be separated according to size with a mass spectrometer that sorts charged particles with approximately the same energy according to their masses. This is usually done by deflecting the charged clusters or ions with an electric or magnetic field; the smaller the mass, the greater is the deflection. This is one of the most effective ways of preparing a beam of clusters of only a single selected mass. It does not eliminate the problem of multiple structures, however.

A technique that can sometimes be used to sort clusters according to their size and structure is a two-step process in which one cluster species at a time is excited with the light from a laser and is then ionized with light from a second laser. This process, called resonant two-photon ionization, is highly selective if the clusters being separated have moderately different absorption spectra. Since this is frequently the case, the method is quite powerful. As the experimenter varies the wavelength of the first exciting laser, a spectrum is produced that includes those wavelengths of light that excite the cluster. If the wavelength of the second ionizing laser is varied, the method also yields the ionization potential, which is the minimum energy that the photon in the ionizing beam must possess in order to knock an electron out of the cluster. Such data help to reveal the forces that bind the cluster together and give some indication of how the cluster will react with atoms, molecules, or other clusters.

A powerful tool for studying clusters is computer simulation of their behaviour. If the nature of the forces between the individual atoms or molecules in a cluster is known, then one can construct a computer model that represents the behaviour of those atoms or molecules by solving the equations of motion of the cluster. To describe the cluster in terms of classical mechanics, the Newtonian equations of motion are solved repeatedly—namely, force equals mass times acceleration, in which the forces depend on the instantaneous positions of all the particles. Hence, these equations are simultaneous, interlinked equations; there is one set of three (for the three instantaneous coordinates of each particle) for each atom or molecule. The results can take one of three forms: (1) the positions and coordinates of the atoms, given in tables, (2) the average properties of the entire cluster, or (3) animations. Tables are too cumbersome for most purposes, and specific average properties are frequently what the investigator seeks. Animated sequences show the same content as the tables but far more efficiently than extensive tables do. In fact, animations sometimes reveal considerably more than is expected by scientists.

It is also possible to construct computer models of clusters based on quantum mechanics instead of Newton's classical mechanics. This is especially appropriate for clusters of hydrogen and helium, because the small masses of their constituent atoms make them very quantumlike in the sense that they reveal the wavelike character that all matter exhibits according to quantum mechanics. The same kinds of data and inferences can be extracted from quantum mechanical calculations as from classical ones, but the preparation and visualization of animations for such clusters are much more demanding than their classical mechanical counterparts.

Structure and properties

STRUCTURE

The abundance distributions for several kinds of clusters show that there are certain sizes of clusters with excep-

tional stability, analogous to the exceptional stability of the atoms of the inert gases helium, neon, argon, krypton, and xenon and of the so-called magic number nuclei—*i.e.*, the sequence of unusually stable atomic nuclei beginning with the α -particle, or helium nucleus. Such unusual stability suggests that its interpretation should be associated with the closing of some kind of shell, or energy level. The overall structure that determines the cluster's stability is generally called its shell structure.

Clusters of atoms bound by van der Waals forces or by other simple forces that depend only on the distance between each pair of atoms have unusual stability when the cluster has exactly the number of atoms needed to form a regular icosahedron. The first three clusters in this series have, respectively, 13, 55, and 147 atoms. These are shown in Figure 47. In the 13-atom cluster, all but one of the atoms occupy equivalent sites. The 55-atom cluster in this series consists of a core—which is just the 13-atom icosahedron—plus 12 more atoms atop the 12 vertices of the icosahedron and 30 more atoms, one in the centre of each of the 30 edges of the icosahedron. The 147-atom cluster consists of a 55-atom icosahedral core, 12 more atoms at the vertices of the outermost shell, one atom in the centre of each of the 20 faces, and two atoms along each of the 30 edges between the vertices. The shell structure that provides special stabilities in this class of clusters is determined by the individual stabilities of the shells of the atoms themselves.

A different kind of extraordinary stability manifests itself in clusters of simple metal atoms. The shell structure for this class of clusters is determined by the electrons and the filling of those shells that have energy states available to the electrons. The numbers of electrons corresponding to closed electron shells in metal clusters are 8, 20, 40, 58, The electron structure can be modeled by supposing that the positively charged cores consisting of the protons and inner-shell electrons of all the cluster's atoms are smeared out into a continuous, attractive background,

Stability of icosahedral structures

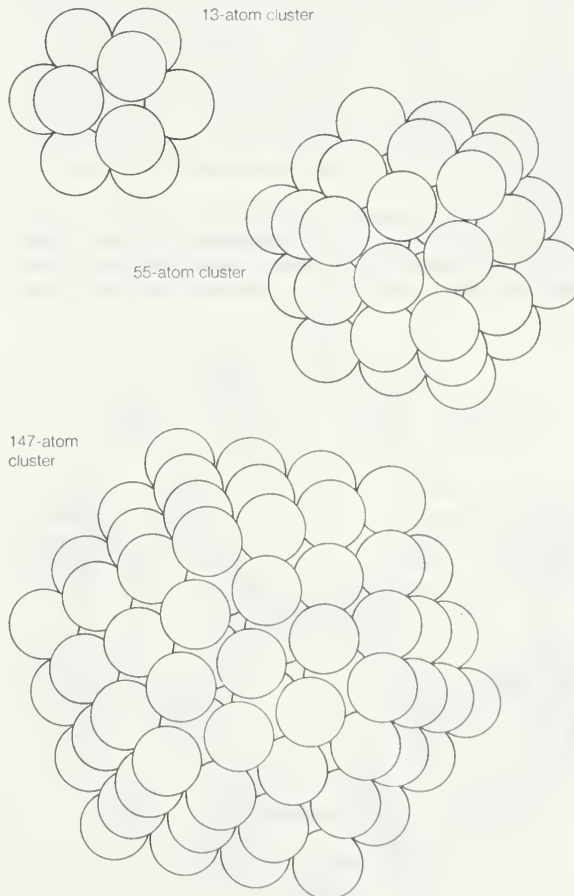


Figure 47: The first three complete icosahedral structures of 13, 55, and 147 particles. These are the structures taken on by clusters of 13, 55, and 147 atoms of neon, argon, krypton, and xenon, for example.

while the valence, or outer-shell, electrons are delocalized (*i.e.*, shared among all atoms in the cluster). The electron environment is much like a well or pit with a flat bottom and a moderately steep wall. The determination of the energy states available for electrons in such a simplified model system is relatively easy and gives a good description of clusters of more than about eight or nine alkali atoms—*i.e.*, lithium, sodium, potassium, rubidium, or cesium. The single valence, or outer-shell, electron of each alkali atom is treated explicitly, while all the others are considered part of the smeared-out core. Since each alkali atom has only one valence electron, the unusually stable clusters of alkalis consist of 8, 20, 40, . . . atoms, corresponding to major shell closings. This model is not as successful in treating metals such as aluminum, which have more than one valence electron.

Still another kind of particularly stable closed shell occurs in clusters sometimes called network structures. The best-known of these is C_{60} , the 60-atom cluster of carbon atoms. In this cluster the atoms occupy the sites of the 60 equivalent vertices of the soccer ball structure, which can be constructed by cutting off the 12 vertices of the icosahedron to make 12 regular 5-sided (regular pentagonal) faces. The icosahedron itself has 20 triangular faces; when its vertices are sliced off, the triangles become hexagons. The 12 pentagons share their edges with these 20 hexagonal faces. No two pentagons have any common edge in this molecule or cluster (C_{60} may be considered either). The resulting high-symmetry structure has been named buckminsterfullerene, after R. Buckminster Fuller, who advocated using such geometric structures in architectural design (see Figure 48).

Other network compounds of carbon are also known. To form a closed-shell structure, a network compound of carbon must have exactly 12 rings of 5 carbon atoms, but the number of rings of 6 carbon atoms is variable. Shells smaller than C_{60} have been discovered, but some of their constituent pentagons must share edges; this makes the smaller network compounds less stable than C_{60} . Shells larger than C_{60} , such as C_{70} , C_{76} , and C_{84} , are known and are relatively stable. Even tubes and "onions" of concentric layers of carbon shells have been reported in observations made with modern electron microscopes known as scanning tunneling microscopes. These devices are powerful enough to reveal images of extremely small clusters and even individual foreign atoms deposited on clean surfaces.

The network compounds of carbon, which make up the class called fullerenes, form compounds with alkali and other metals. Some of these compounds of fullerenes combined with metals, such as K_3C_{60} , become superconductors

Fullerenes

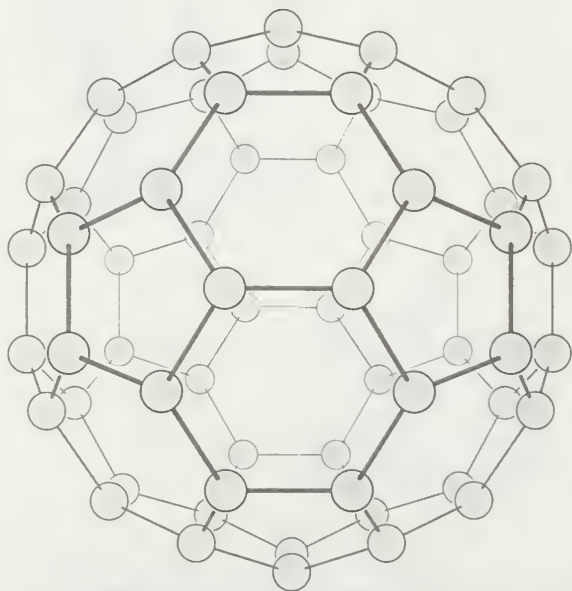


Figure 48: The structure of C_{60} , buckminsterfullerene. The geometry is that of a soccer ball with a carbon atom at each vertex (see text).

at low temperatures; that is to say, they lose all resistance to electric current flow when they are cooled sufficiently. The class of network compounds as a group had been imagined from time to time, but only in the late 1980s were they realized in the laboratory and shown to have closed-shell network structures.

PHYSICAL PROPERTIES

Clusters share some of the physical properties of bulk matter, a few of which are rather surprising. Clusters of all substances except helium and possibly hydrogen are solidlike at low temperatures as expected. The atoms or molecules of a cluster remain close to their equilibrium positions, vibrating around these positions in moderately regular motions of small amplitude. This is characteristic of all solids; their atoms are constrained to stay roughly in the same position at all times. In a liquid or a gas, the atoms or molecules are free to wander through the space accessible to the substance. A gas or vapour has so much empty space relative to the volume occupied by the particles that the particles move almost unhindered, colliding only occasionally with other particles or with the walls of the container. A liquid is typically almost as dense as a solid but has some empty spaces into which the atoms or molecules can easily move. Hence, the particles of a liquid can diffuse with moderate ease. (Water is an exception; its density as a liquid is higher than its density as ice, because ice has an unusually open structure in comparison with most solids, and this open structure collapses when ice turns to water.) Clusters can be liquidlike if they are warm enough, but typically the temperatures at which clusters can become liquid are much lower than the melting points of the corresponding bulk solids. If temperatures are measured on the Kelvin scale, small clusters become liquidlike at temperatures of roughly half the bulk melting temperatures. For example, solid argon melts at approximately 80 K, while small clusters of argon become liquid at about 40 K.

Some clusters are expected to show a gradual transition from solidlike to liquidlike, appearing slushy in the temperature range between their solidlike and liquidlike zones. Other clusters are expected to show, as seen in computer simulations, distinct solidlike and liquidlike forms that qualitatively resemble bulk solids and liquids in virtually every aspect, even though they may exhibit quantitative differences from the bulk. Solid clusters, for example, show virtually no diffusion, but the particles of a liquid cluster can and do diffuse. The forces that hold a particle in place in a solid cluster are strong, comparable to those of a bulk solid; but those in a liquid cluster include, in addition to forces comparable in strength to those in solids, some forces weak enough to allow a particle to wander far from its home base and find new equilibrium positions. Those same weak forces are responsible for making a liquid cluster compliant; that is, weak forces allow the liquid to accommodate any new force, say, a finger inserted into water. Ice will not yield to such an intruding force, but when a finger is placed into liquid water, the water molecules move aside under the force of the finger. This is much like the behaviour of a bulk liquid. The greatest differences between bulk solids and liquids and solid and liquid clusters arise from the fact that a large fraction of the particles of a cluster are on its surface. As a result, the particle mobility that characterizes liquids and enables them to exhibit diffusion and physical compliance is enhanced in a cluster, for the cluster can easily expand by enlarging the spaces between particles and can also transfer particles from its interior to its surface, leaving vacancies that enhance the mobility of the interior particles. The large surface area, together with the curved shape of the cluster's surface, make it easier for particles to leave a cluster than to leave the flat surface of a bulk liquid or solid. An important consequence is that the vapour pressure of a cluster is higher than the vapour pressure of the corresponding bulk, and accordingly the boiling point of a liquid cluster—*i.e.*, the temperature at which the vapour pressure of a liquid is equal to the pressure of the surrounding atmosphere—is lower than that of the corresponding bulk liquid. The vapour pressure of

clusters decreases with increasing cluster size, while the boiling point increases.

Perhaps the greatest difference between clusters and bulk matter with regard to their transformation between solid and liquid is the nature of the equilibrium between two phases. Bulk solids can be in equilibrium with their liquid forms at only a single temperature for any given pressure or at only a single pressure for any given temperature. A graph of the temperatures and pressures along which the solid and liquid forms of any given substance are in equilibrium is called a coexistence curve. One point on the coexistence curve for ice and liquid water is 0° C and one atmosphere of pressure. A similar curve can be drawn for the coexistence of any two bulk phases, such as liquid and vapour; a point on the coexistence curve for liquid water and steam is 100° C and one atmosphere of pressure. Clusters differ sharply from bulk matter in that solid and liquid clusters of the same composition are capable of coexisting within a band of temperatures and pressures. At any chosen pressure, the proportion of liquid clusters to solid clusters increases with temperature. At low temperatures the clusters are solid, as described above. As the temperature is increased, some clusters transform from solid to liquid. If the temperature is raised further, the proportion of liquid clusters increases, passing through 50 percent, so that the mixture becomes predominantly liquid clusters. At sufficiently high temperatures all the clusters are liquid.

No cluster remains solid or liquid all the time; liquidlike clusters occasionally transform spontaneously into solidlike clusters and vice versa. The fraction of time that a particular cluster spends as a liquid is precisely the same as the fraction of clusters of that same type within a large collection that are liquid at a given instant. That is to say, the time average behaviour gives the same result as the ensemble average, which is the average over a large collection of identical objects. This equivalence is not limited to clusters; it is the well-known ergodic property that is expected of all but the simplest real systems.

Other significant physical properties of clusters are their electric, magnetic, and optical properties. The electric properties of clusters, such as their conductivity and metallic or insulating character, depend on the substance and the size of the cluster. Quantum theory attributes wavelike character to matter, a behaviour that is detectable only when matter is examined on the scale of atoms and electrons. At a scale of millimetres or even millionths of millimetres, the wavelengths of matter are too short to be observed. Clusters are often much smaller than that, with the important consequence that many are so small that when examined their electrons and electronic states can exhibit the wavelike properties of matter. In fact, quantum properties may play an important role in determining the electrical character of the cluster. In particular, as described previously, if a cluster is extremely small, the energy levels or quantum states of its electrons are not close enough together to permit the cluster to conduct electricity.

Moreover, an alternative way to view this situation is to recognize that a constant electric force (*i.e.*, the kind that drives a direct current) and an alternating force (the kind that generates alternating current) can behave quite differently in a cluster. Direct current cannot flow in an isolated cluster and probably cannot occur in a small cluster even if it is sandwiched between slabs of metal. The current flow is prohibited both because the electrons that carry the current encounter the boundaries of the cluster and because there are no quantum states readily available at energies just above those of the occupied states, which are the states that must be achieved to allow the electrons to move. However, if a field of alternating electric force is applied with a frequency of alternation so high that the electrons are made to reverse their paths before they encounter the boundaries of the cluster, then the equivalent of conduction will take place. Ordinary 60-cycle (60-hertz) alternating voltage and even alternations at radio-wave frequencies switch direction far too slowly to produce this behaviour in clusters; microwave frequencies are required.

Magnetic properties of clusters, in contrast, appear to be rather similar to those of bulk matter. They are not

identical, because clusters contain only small numbers of electrons, which are the particles whose magnetic character makes clusters and bulk matter magnetic. As a result, the differences between magnetic properties of clusters and of bulk matter are more a matter of degree than of kind. Clusters of substances magnetic in the bulk also tend to be magnetic. The larger the cluster, the more nearly will the magnetic character per atom approach that of the bulk. The degree of this magnetic character depends on how strongly the individual electron magnets couple to each other to become aligned in the same direction; the larger the cluster, the stronger is this coupling.

The optical properties of weakly bound clusters are much like those of their component atoms or molecules; the small differences are frequently useful diagnostics of how the cluster is bound and what its structure may be. Optical properties of metal clusters are more like those of the corresponding bulk metals than like those of the constituent atoms. These properties reveal which cluster sizes are unusually stable and therefore correspond to "magic-number" sizes. Optical properties of covalently bound clusters are in most cases—*e.g.*, fullerenes—unlike those of either the component atoms or the bulk but are important clues to the structure and bonding of the cluster.

CHEMICAL PROPERTIES

The chemical properties of clusters are a combination of the properties of bulk and molecular matter. Several kinds of clusters, particularly those of the metallic variety, induce certain molecules to dissociate. For example, hydrogen molecules, H₂, spontaneously break into two hydrogen atoms when they attach themselves to a cluster of iron atoms. Ammonia likewise dissociates when attaching itself to an iron cluster. Similar reactions occur with bulk matter, but the rate at which such gases react with bulk metals depends only on how much gaseous reactant is present and how much surface area the bulk metal presents to the gas. Metal powders react much faster than dense solids with the same total mass because they have much more surface area. Small and medium-size clusters, on the other hand, show different reactivities for every size, although these reactivities do not vary smoothly with size. Furthermore, there are instances, such as reactions of hydrogen with iron, in which two different geometric forms of clusters of a single size have different reaction rates, just as two different molecules with the same elemental composition, called chemical isomers, may have different reaction rates with the same reactant partner. In the case of molecules, this is not surprising, because different isomers typically have quite different structures, physical properties, and reactivities and do not normally transform readily into one another. Isomers of clusters of a specific chemical composition, however, may well transform into one another with moderate ease and with no excessive increase in energy above the amount present when they formed. If the reaction releases energy (*i.e.*, it is exothermic) in sufficient quantity to transform the cluster from solid to liquid, a cluster may melt as it reacts.

Some of the interesting chemistry of clusters is set in motion by light. For example, light of sufficiently short wavelength can dissociate molecules that are captured in the middle of a cluster of nonreactive atoms or molecules. A common question is whether the surrounding molecules or atoms form a cage strong enough to prevent the fragment atoms from flying apart and from leaving the cluster. The answer is that, if there are only a few surrounding atoms or molecules, the fragments escape their initial cage, and, if the energy of the light is high enough, at least one of the fragments escapes. On the other hand, if there are enough nonreacting cage atoms or molecules in the cluster to form at least one complete shell around the molecule that breaks up, the cage usually holds the fragments close together until they eventually recombine.

A related sort of reaction, another example of competition between a particle's attempt to escape from a cluster and some other process, occurs if light is used to detach an excess electron from a negative ion in the middle of an inert cluster. If, for example, light knocks the extra electron off a free, negatively charged bromine molecule, Br₂⁻, the

Optical properties

Cluster-induced dissociation

Electric properties

electron of course escapes. If the charged molecule is surrounded by a few inert molecules of carbon dioxide (CO_2), the electron escapes almost as readily. If 10 or 15 CO_2 molecules encase the Br_2^- , the electron does not escape; instead, it loses its energy to the surrounding molecules of CO_2 , some of which boil off, and then eventually recombines with the now neutral bromine molecule to re-form the original Br_2^- .

The chemical properties of fullerenes and other network compounds have become a subject of their own, bridging molecular and cluster chemistry. These compounds typically react with a specific number of other atoms or molecules to form new species with definite compositions and structures. Compounds such as K_3C_{60} mentioned previously have the three potassium atoms outside the C_{60} cage, all as singly charged ions, K^+ , and the ball of 60 carbon atoms carries three negative charges to make the entire compound electrically neutral. Other compounds of C_{60} , such as that made with the metal lanthanum, contain the metal inside the carbon cage, forming a new kind of substance. It is possible to add or take away hydrogen atoms from C_{60} and its larger relatives, much as hydrogen atoms can be added or removed from some kinds of hydrocarbons; in this way some of the chemistry of this class of clusters is similar to classical organic chemistry.

One of the goals of cluster science is the creation of new

kinds of materials. The possible preparation of diamond films is one such application; another example is the proposal to make so-called superatoms that consist of an electron donor atom in the centre of a cluster of electron acceptors; the fullerene clusters containing a metal ion inside the cage seem to be just such a species but with much more open structures than had been previously envisioned. Molecular electronics is another goal; in this technology clusters would be constructed with electrical properties much like those of transistors and could be packed together to make microcircuits far smaller than any now produced. These applications are still theoretical, however, and have not yet been realized.

Clusters do indeed form a bridge between bulk and molecular matter. Their physical and chemical properties are in many instances unique to their finely divided state. Some examples of clusters, such as the network clusters of carbon, are new forms of matter. Nevertheless, such clusters, particularly the small and middle-size ones, not only exhibit behaviours of their own but also provide new insights into the molecular origins of the properties of bulk matter. They may yield other new materials—*e.g.*, possibly far more disordered, amorphous glasslike substances than the glasses now in common use—and at the same time give rise to deeper understanding of why and how glasses form at all. (R.S.Be.)

LOW-TEMPERATURE PHENOMENA

The term low-temperature phenomena refers to the behaviour of matter at temperatures closer to absolute zero (-273.15°C [-459.67°F]) than to room temperature. At such temperatures the thermal, electric, and magnetic properties of many substances undergo great change, and, indeed, the behaviour of matter may seem strange when compared with that at room temperature. Superconductivity and superfluidity can be cited as two such phenomena that occur below certain critical temperatures; in the former, many chemical elements, compounds, and alloys show no resistance whatsoever to the flow of electricity, and, in the latter, liquid helium can flow through tiny holes impervious to any other liquid.

Superconductivity

Superconductivity was discovered in 1911 by the Dutch physicist Heike Kamerlingh Onnes; he was awarded the Nobel Prize for Physics in 1913 for his low-temperature research. Kamerlingh Onnes found that the electrical resistivity of a mercury wire disappears suddenly when it is cooled below a temperature of about 4 K (-269°C); absolute zero is 0 K, the temperature at which all matter loses its disorder. He soon discovered that a superconducting material can be returned to the normal (*i.e.*, nonsuperconducting) state either by passing a sufficiently large current through it or by applying a sufficiently strong magnetic field to it.

For many years it was believed that, except for the fact that they had no electrical resistance (*i.e.*, that they had infinite electrical conductivity), superconductors had the same properties as normal materials. This belief was shattered in 1933 by the discovery that a superconductor is highly diamagnetic; that is, it is strongly repelled by and tends to expel a magnetic field. This phenomenon, which is very strong in superconductors, is called the Meissner effect for one of the two men who discovered it. Its discovery made it possible to formulate, in 1934, a theory of the electromagnetic properties of superconductors that predicted the existence of an electromagnetic penetration depth (see below *The Meissner effect*), which was first confirmed experimentally in 1939. In 1950 it was clearly shown for the first time that a theory of superconductivity must take into account the fact that free electrons in a crystal are influenced by the vibrations of atoms that define the crystal structure, called the lattice vibrations. In 1953, in an analysis of the thermal conductivity of superconductors, it was recognized that the distribution of

energies of the free electrons in a superconductor is not uniform but has a separation called the energy gap.

The theories referred to thus far served to show some of the interrelationships between observed phenomena but did not explain them as consequences of the fundamental laws of physics. For almost 50 years after Kamerlingh Onnes' discovery, theorists were unable to develop a fundamental theory of superconductivity. Finally, in 1957 such a theory was presented by the physicists John Bardeen, Leon N. Cooper, and John Robert Schrieffer of the United States; it won for them the Nobel Prize for Physics in 1972. It is now called the BCS theory in their honour, and most later theoretical work is based on it. The BCS theory also provided a foundation for an earlier model that had been introduced by the Russian physicists Lev Davidovich Landau and Vitaly Lazarevich Ginzburg (1950). This model has been useful in understanding electromagnetic properties, including the fact that any internal magnetic flux in superconductors exists only in discrete amounts (instead of in a continuous spectrum of values), an effect called the quantization of magnetic flux. This flux quantization, which had been predicted from quantum mechanical principles, was first observed experimentally in 1961.

In 1962 the British physicist Brian D. Josephson predicted that two superconducting objects placed in electric contact would display certain remarkable electromagnetic properties. These properties have since been observed in a wide variety of experiments, demonstrating quantum mechanical effects on a macroscopic scale.

The theory of superconductivity has been tested in a wide range of experiments, involving, for example, ultrasonic absorption studies, nuclear-spin phenomena, low-frequency infrared absorption, and electron-tunneling experiments (see below *Energy gaps*). The results of these measurements have brought understanding to many of the detailed properties of various superconductors.

THERMAL PROPERTIES OF SUPERCONDUCTORS

Superconductivity is a startling departure from the properties of normal (*i.e.*, nonsuperconducting) conductors of electricity. In materials that are electric conductors, some of the electrons are not bound to individual atoms but are free to move through the material; their motion constitutes an electric current. In normal conductors these so-called conduction electrons are scattered by impurities, dislocations, grain boundaries, and lattice vibrations (phonons). In a superconductor, however, there is an ordering among

The BCS theory

the conduction electrons that prevents this scattering. Consequently, electric current can flow with no resistance at all. The ordering of the electrons, called Cooper pairing, involves the momenta of the electrons rather than their positions. The energy per electron that is associated with this ordering is extremely small, typically about one thousandth of the amount by which the energy per electron changes when a chemical reaction takes place. One reason that superconductivity remained unexplained for so long is the smallness of the energy changes that accompany the transition between normal and superconducting states. In fact, many incorrect theories of superconductivity were advanced before the BCS theory was proposed. For additional details on electric conduction in metals and the effects of temperature and other influences, see the article **ELECTRICITY AND MAGNETISM**.

Superconducting materials

Hundreds of materials are known to become superconducting at low temperatures. Twenty-seven of the chemical elements, all of them metals, are superconductors in their usual crystallographic forms at low temperatures and low (atmospheric) pressure. Among these are commonly known metals such as aluminum, tin, lead, and mercury and less common ones such as rhenium, lanthanum, and protactinium. In addition, 11 chemical elements that are metals, semimetals, or semiconductors are superconductors at low temperatures and high pressures. Among these are uranium, cerium, silicon, and selenium. Bismuth and five other elements, though not superconducting in their usual crystallographic form, can be made superconducting by preparing them in a highly disordered form, which is stable at extremely low temperatures. Superconductivity is not exhibited by any of the magnetic elements chromium, manganese, iron, cobalt, or nickel.

Most of the known superconductors are alloys or compounds. It is possible for a compound to be superconducting even if the chemical elements constituting it are not; examples are disilver fluoride (Ag_2F) and a compound of carbon and potassium (C_8K). Some semiconducting compounds, such as tin telluride (SnTe), become superconducting if they are properly doped with impurities.

Since 1986 some compounds containing copper and oxygen (called cuprates) have been found to have extraordinarily high transition temperatures, denoted T_c . This is the temperature below which a substance is superconducting. The properties of these high- T_c compounds are different in some respects from those of the types of superconductors known prior to 1986, which will be referred to as classic superconductors in this discussion. For the most part, the high- T_c superconductors are treated explicitly toward the end of this section. In the discussion that immediately follows, the properties possessed by both kinds of superconductors will be described, with attention paid to specific differences for the high- T_c materials. A further classification problem is presented by the superconducting compounds of carbon (sometimes doped with other atoms) in which the carbon atoms are on the surface of a cluster with a spherical or spheroidal crystallographic structure. These compounds, discovered in the 1980s, are called fullerenes (if only carbon is present) or fullerides (if doped). They have superconducting transition temperatures higher than those of the classic superconductors. It is not yet known whether these compounds are fundamentally similar to the cuprate high-temperature superconductors.

Transition temperatures. The vast majority of the known superconductors have transition temperatures that lie between 1 K and 10 K. Of the chemical elements,

tungsten has the lowest transition temperature, 0.015 K, and niobium the highest, 9.2 K. The transition temperatures of some of the commonly known superconducting elements are indicated in Table 9.

The transition temperature is usually very sensitive to the presence of magnetic impurities. A few parts per million of manganese in zinc, for example, lowers the transition temperature considerably.

Specific heat and thermal conductivity. The thermal properties of a superconductor can be compared with those of the same material at the same temperature in the normal state. (The material can be forced into the normal state at low temperature by a large enough magnetic field.)

Adapted from *Zeitschrift für Physik* (1959)

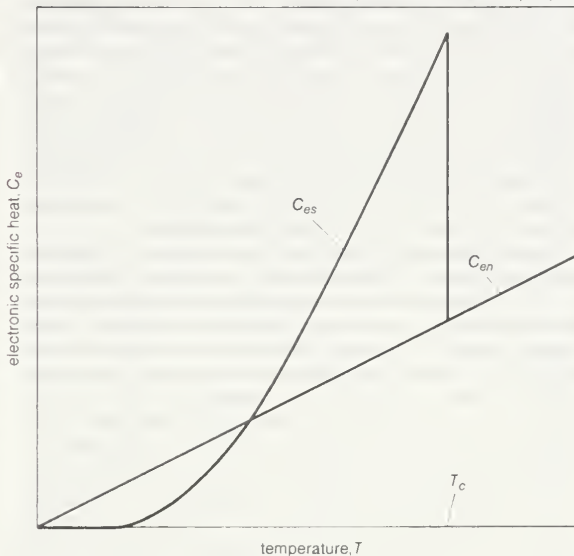


Figure 49: Specific heat in the normal (C_{en}) and superconducting (C_{es}) states of a classic superconductor as a function of absolute temperature. The two functions are identical at the transition temperature (T_c) and above T_c .

When a small amount of heat is put into a system, some of the energy is used to increase the lattice vibrations (an amount that is the same for a system in the normal and in the superconducting state), and the remainder is used to increase the energy of the conduction electrons. The electronic specific heat (C_e) of the electrons is defined as the ratio of that portion of the heat used by the electrons to the rise in temperature of the system. Figure 49 shows how the specific heat of the electrons in a superconductor varies with the absolute temperature (T) in the normal and in the superconducting state. It is evident from the figure that the electronic specific heat in the superconducting state (designated C_{es}) is smaller than in the normal state (designated C_{en}) at low enough temperatures but that C_{es} becomes larger than C_{en} as the transition temperature T_c is approached, at which point it drops abruptly to C_{en} for the classic superconductors, although the curve has a cusp shape near T_c for the high- T_c superconductors. Precise measurements have indicated that, at temperatures considerably below the transition temperature, the logarithm of the electronic specific heat is inversely proportional to the temperature. This temperature dependence, together with the principles of statistical mechanics, strongly suggests that there is a gap in the distribution of energy levels available to the electrons in a superconductor, so that a minimum energy is required for the excitation of each electron from a state below the gap to a state above the gap. Some of the high- T_c superconductors provide an additional contribution to the specific heat, which is proportional to the temperature. This behaviour indicates that there are electronic states lying at low energy; additional evidence of such states is obtained from optical properties and tunneling measurements.

The heat flow per unit area of a sample equals the product of the thermal conductivity (K) and the temperature gradient ∇T : $J_Q = -K \nabla T$, the minus sign indicating that heat always flows from a warmer to a colder region of a substance.

Temperature dependence of electronic specific heat

Table 9: Transition Temperatures and Low-Temperature Values of Critical Magnetic Fields of Some Superconducting Elements

	T_c (K)	H_0 (oersted)
Zinc	0.88	53
Aluminum	1.20	99
Indium	3.41	282
Tin	3.72	306
Mercury	4.15	411
Lead	7.19	803

The thermal conductivity in the normal state (K_n) approaches the thermal conductivity in the superconducting state (K_s) as the temperature (T) approaches the transition temperature (T_c) for all materials, whether they are pure or impure. This suggests that the energy gap (Δ) for each electron approaches zero as the temperature (T) approaches the transition temperature (T_c). This would also account for the fact that the electronic specific heat in the superconducting state (C_{en}) is higher than in the normal state (C_{en}) near the transition temperature: as the temperature is raised toward the transition temperature (T_c), the energy gap in the superconducting state decreases, the number of thermally excited electrons increases, and this requires the absorption of heat.

Energy gaps. As stated above, the thermal properties of superconductors indicate that there is a gap in the distribution of energy levels available to the electrons, and so a finite amount of energy, designated as delta (Δ), must be supplied to an electron to excite it. This energy is maximum (designated Δ_0) at absolute zero and changes little with increase of temperature until the transition temperature is approached, where Δ decreases to zero, its value in the normal state. The BCS theory predicts an energy gap with just this type of temperature dependence.

According to the BCS theory, there is a type of electron pairing (electrons of opposite spin acting in unison) in the superconductor that is important in interpreting many superconducting phenomena. The electron pairs, called Cooper pairs, are broken up as the superconductor is heated. Each time a pair is broken, an amount of energy that is at least as much as the energy gap (Δ) must be supplied to each of the two electrons in the pair, so an energy at least twice as great (2Δ) must be supplied to the superconductor. The value of twice the energy gap at 0 K (which is $2\Delta_0$) might be assumed to be higher when the transition temperature of the superconductor is higher. In fact, the BCS theory predicts a relation of this type—namely, that the energy supplied to the superconductor at absolute zero would be $2\Delta_0 = 3.53 kT_c$, where k is Boltzmann's constant (1.38×10^{-23} joule per kelvin). In the high- T_c cuprate compounds, values of $2\Delta_0$ range from approximately three to eight multiplied by kT_c .

The energy gap (Δ) can be measured most precisely in a tunneling experiment (a process in quantum mechanics that allows an electron to escape from a metal without acquiring the energy required along the way according to the laws of classical physics). In this experiment, a thin insulating junction is prepared between a superconductor and another metal, assumed here to be in the normal state. In this situation, electrons can quantum mechanically tunnel from the normal metal to the superconductor if they have sufficient energy. This energy can be supplied by applying a negative voltage (V) to the normal metal, with respect to the voltage of the superconductor.

Tunneling will occur if eV —the product of the electron charge, e (-1.60×10^{-19} coulomb), and the voltage—is at least as large as the energy gap Δ . The current flowing between the two sides of the junction is small up to a voltage equal to $V = \Delta/e$, but then it rises sharply. This provides an experimental determination of the energy gap (Δ). In describing this experiment it is assumed here that the tunneling electrons must get their energy from the applied voltage rather than from thermal excitation.

MAGNETIC AND ELECTROMAGNETIC PROPERTIES OF SUPERCONDUCTORS

Critical field. One of the ways in which a superconductor can be forced into the normal state is by applying a magnetic field. The weakest magnetic field that will cause this transition is called the critical field (H_c) if the sample is in the form of a long, thin cylinder or ellipsoid and the field is oriented parallel to the long axis of the sample. (In other configurations the sample goes from the superconducting state into an intermediate state, in which some regions are normal and others are superconducting, and finally into the normal state.) The critical field increases with decreasing temperature. For the superconducting elements, its values (H_0) at absolute zero range from 1.1 oersted for tungsten to 830 oersteds for tantalum. Values

of H_0 are listed in Table 9 for some common superconducting elements.

These remarks about the critical field apply to ordinary (so-called type I) superconductors. In the following section the behaviour of other (type II) superconductors is examined.

The Meissner effect. As was stated above, a type I superconductor in the form of a long, thin cylinder or ellipsoid remains superconducting at a fixed temperature as an axially oriented magnetic field is applied, provided the applied field does not exceed a critical value (H_c). Under these conditions, superconductors exclude the magnetic field from their interior, as could be predicted from the laws of electromagnetism and the fact that the superconductor has no electric resistance. A more astonishing effect occurs if the magnetic field is applied in the same way to the same type of sample at a temperature above the transition temperature and is then held at a fixed value while the sample is cooled. It is found that the sample expels the magnetic flux as it becomes superconducting. This is called the Meissner effect. Complete expulsion of the magnetic flux (a complete Meissner effect) occurs in this way for certain superconductors, called type I superconductors, but only for samples that have the described geometry. For samples of other shapes, including hollow structures, some of the magnetic flux can be trapped, producing an incomplete or partial Meissner effect.

Type II superconductors have a different magnetic behaviour. Examples of materials of this type are niobium and vanadium (the only type II superconductors among the chemical elements) and some alloys and compounds, including the high- T_c compounds. As a sample of this type, in the form of a long, thin cylinder or ellipsoid, is exposed to a decreasing magnetic field that is axially oriented with the sample, the increase of magnetization, instead of occurring suddenly at the critical field (H_c), sets in gradually. Beginning at the upper critical field (H_{c2}), it is completed at a lower critical field (H_{c1} ; see Figure 50). If the sample is of some other shape, is hollow, or is inhomogeneous or strained, some magnetic flux remains trapped, and some magnetization of the sample remains after the applied field is completely removed. Known values of the upper critical field extend up to 6×10^5 oersteds, the value for the compound of lead, molybdenum, and sulfur with formula $PbMo_6S_8$.

Cooper
electron
pairs

Behaviour
of type II
super-
conductors

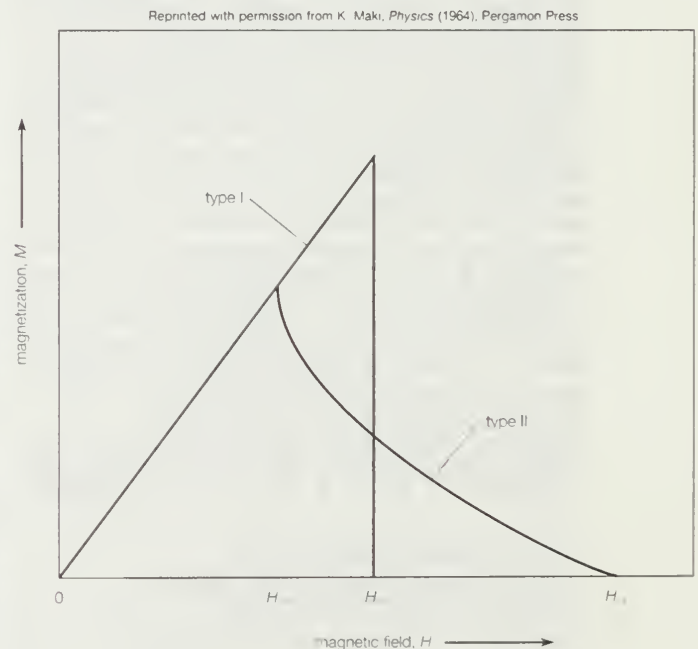


Figure 50: Magnetization as a function of magnetic field for a type I superconductor and a type II superconductor. For type I, magnetic flux is expelled, producing a magnetization (M) that increases with magnetic field (H) until a critical field (H_c) is reached, at which it falls to zero as with a normal conductor. A type II superconductor has two critical magnetic fields (H_{c1} and H_{c2}); below H_{c1} type II behaves as type I, and above H_{c2} it becomes normal.

Electro-
magnetic
penetration
depth

The expulsion of magnetic flux by type I superconductors in fields below the critical field (H_c) or by type II superconductors in fields below H_{c1} is never quite as complete as has been stated in this simplified presentation, because the field always penetrates into a sample for a small distance, known as the electromagnetic penetration depth. Values of the penetration depth for the superconducting elements at low temperature lie in the range from about 390 to 1300 angstroms. As the temperature approaches the critical temperature, the penetration depth becomes extremely large.

High-frequency electromagnetic properties. The foregoing descriptions have pertained to the behaviour of superconductors in the absence of electromagnetic fields or in the presence of steady or slowly varying fields; the properties of superconductors in the presence of high-frequency electromagnetic fields, however, have also been studied.

The energy gap in a superconductor has a direct effect on the absorption of electromagnetic radiation. At low temperatures, at which a negligible fraction of the electrons are thermally excited to states above the gap, the superconductor can absorb energy only in a quantized amount that is at least twice the gap energy (at absolute zero, $2\Delta_0$). In the absorption process, a photon (a quantum of electromagnetic energy) is absorbed, and a Cooper pair is broken; both electrons in the pair become excited. The photon's energy (E) is related to its frequency (ν) by the Planck relation, $E = h\nu$, in which h is Planck's constant (6.63×10^{-34} joule second). Hence the superconductor can absorb electromagnetic energy only for frequencies at least as large as $2\Delta_0/h$.

Magnetic-flux quantization. The laws of quantum mechanics dictate that electrons have wave properties and that the properties of an electron can be summed up in what is called a wave function. If several wave functions are in phase (*i.e.*, act in unison), they are said to be coherent. The theory of superconductivity indicates that there is a single, coherent, quantum mechanical wave function that determines the behaviour of all the superconducting electrons. As a consequence, a direct relationship can be shown to exist between the velocity of these electrons and the magnetic flux (Φ) enclosed within any closed path inside the superconductor. Indeed, inasmuch as the magnetic flux arises because of the motion of the electrons, the magnetic flux can be shown to be quantized; *i.e.*, the intensity of this trapped flux can change only by units of Planck's constant divided by twice the electron charge.

When a magnetic field enters a type II superconductor (in an applied field between the lower and upper critical fields, H_{c1} and H_{c2}), it does so in the form of quantized fluxoids, each carrying one quantum of flux. These fluxoids tend to arrange themselves in regular patterns that have been detected by electron microscopy and by neutron diffraction. If a large enough current is passed through the superconductor, the fluxoids move. This motion leads to energy dissipation that can heat the superconductor and drive it into the normal state. The maximum current per unit area that a superconductor can carry without being forced into the normal state is called the critical current density (J_c). In making wire for superconducting high-field magnets, manufacturers try to fix the positions of the fluxoids by making the wire inhomogeneous in composition.

Josephson currents. If two superconductors are separated by an insulating film that forms a low-resistance junction between them, it is found that Cooper pairs can tunnel from one side of the junction to the other. (This process occurs in addition to the single-particle tunneling already described.) Thus, a flow of electrons, called the Josephson current, is generated and is intimately related to the phases of the coherent quantum mechanical wave function for all the superconducting electrons on the two sides of the junction. It was predicted that several novel phenomena should be observable, and experiments have demonstrated them. These are collectively called the Josephson effect or effects.

The first of these phenomena is the passage of current through the junction in the absence of a voltage across the junction. The maximum current that can flow at zero voltage depends on the magnetic flux (Φ) passing through

the junction as a result of the magnetic field generated by currents in the junction and elsewhere. The dependence of the maximum zero-voltage current on the magnetic field applied to a junction between two superconductors is shown in Figure 51.

A second type of Josephson effect is an oscillating current resulting from a relation between the voltage across the junction and the frequency (ν) of the currents associated with Cooper pairs passing through the junction. The frequency (ν) of this Josephson current is given by $\nu = 2eV/h$, where e is the charge of the electron. Thus, the frequency increases by 4.84×10^{14} hertz (cycles per second) for each additional volt applied to the junction. This effect can be demonstrated in various ways. The voltage can be established with a source of direct-current (DC) power, for instance, and the oscillating current can be detected by the electromagnetic radiation of frequency (ν) that it generates. Another method is to expose the junction to radiation of another frequency (ν') generated externally. It is found that a graph of the DC current versus voltage has current steps at values of the voltage corresponding to Josephson frequencies that are integral multiples (n) of the external frequency ($\nu = n\nu'$); that is, $V = nh\nu'/2e$. The observation of current steps of this type has made it possible to measure h/e with far greater precision than by any other method and has therefore contributed to a knowledge of the fundamental constants of nature.

The Josephson effect has been used in the invention of novel devices for extremely high-sensitivity measurements of currents, voltages, and magnetic fields.

Oscillating
currents

Fluxoids



Figure 51: Maximum zero-voltage (Josephson) current passing through a junction by Cooper-pair tunneling as a function of magnetic field.

HIGHER-TEMPERATURE SUPERCONDUCTIVITY

Ever since Kamerlingh Onnes discovered that mercury becomes superconducting at temperatures less than 4 K, scientists have been searching for superconducting materials with higher transition temperatures. Until 1986 a compound of niobium and germanium (Nb_3Ge) had the highest known transition temperature, 23 K, less than a 20-degree increase in 75 years. Most researchers expected that the next increase in transition temperature would be found in a similar metallic alloy and that the rise would be only one or two degrees. In 1986, however, the Swiss physicist Karl Alex Müller and his West German associate, Johannes Georg Bednorz, discovered, after a three-year search among metal oxides, a material that had an unprecedentedly high transition temperature of about 30 K. They were awarded the Nobel Prize for Physics in 1987, and their discovery immediately stimulated groups of investigators in China, Japan, and the United States to produce superconducting oxides with even higher transition temperatures.

These high-temperature superconductors are ceramics. They contain lanthanum, yttrium, or another of the rare-earth elements or bismuth or thallium; usually barium or strontium (both alkaline-earth elements); copper; and oxygen. Other atomic species can sometimes be introduced by chemical substitution while retaining the high- T_c properties. The superconducting transition temperatures of some of the high- T_c materials are listed in Table 10. The value 134 K is the highest known T_c value. The compounds given in the table are members of the major families of high- T_c superconductors. Within each family, only the subscripts (*i.e.*, stoichiometry) vary from one compound to another. The compounds listed have the highest observed superconducting transition temperature in their respective families. Samples in the families containing bismuth or thallium always exhibit a great deal of atomic disorder, with atoms in the "wrong" crystallographic sites and with impurity phases. It is possible that such disorder is required to make these compounds thermodynamically stable.

Table 10: Transition Temperatures of Some High- T_c Superconductors

compound	T_c (K)
$\text{Nd}_{1.85}\text{Ce}_{0.15}\text{CuO}_4$	24
$\text{La}_{1.85}\text{Sr}_{0.15}\text{CuO}_4$	40
$\text{YBa}_2\text{Cu}_3\text{O}_7$	92
$\text{Bi}_2\text{Sr}_2\text{Ca}_2\text{Cu}_3\text{O}_{10}$	110
$\text{Tl}_2\text{Ba}_2\text{Ca}_2\text{Cu}_3\text{O}_{10}$	127
$\text{Hg}_2\text{Ba}_2\text{Ca}_2\text{Cu}_3\text{O}_8$	134

The compounds have crystal structures containing planes of Cu and O atoms, and some also have chains of Cu and O atoms. The roles played by these planes and chains have come under intense investigation. Varying the oxygen content or the heat treatment of the materials dramatically changes their transition temperatures, critical magnetic fields, and other properties. Single crystals of the high-temperature superconductors are very anisotropic—*i.e.*, their properties associated with a direction, such as the critical fields or the critical current density, are highly dependent on the angle between that direction and the rows of atoms in the crystal.

If the number of superconducting electrons per unit volume is locally disturbed by an applied force (typically electric or magnetic), this disturbance propagates for a certain distance in the material; the distance is called the superconducting coherence length (or Ginzburg-Landau coherence length), ξ . If a material has a superconducting region and a normal region, many of the superconducting properties disappear gradually—over a distance ξ —upon traveling from the former to the latter region. In the pure (*i.e.*, undoped) classic superconductors ξ is on the order of a few thousand angstroms, but in the high- T_c superconductors it is on the order of 1 to 10 angstroms. The small size of ξ affects the thermodynamic and electromagnetic properties of the high- T_c superconductors. For example, it is responsible for the cusp shape of the specific heat curve near T_c that was mentioned above. It is also responsible for the ability of the high- T_c superconductors to remain superconducting in extraordinarily large fields—on the order of 1,000,000 gauss (100 teslas)—at low temperatures.

The high- T_c superconductors are type II superconductors. They exhibit zero resistance, strong diamagnetism, the Meissner effect, magnetic flux quantization, the Josephson effects, an electromagnetic penetration depth, an energy gap for the superconducting electrons, and the characteristic temperature dependences of the specific heat and the thermal conductivity that are described above. Therefore, it is clear that the conduction electrons in these materials form the Cooper pairs used to explain superconductivity in the BCS theory. Thus, the central conclusions of the BCS theory are demonstrated. Indeed, that theory guided Bednorz and Müller in their search for high-temperature superconductors. It is not known, however, why the transition temperatures of these oxides are so high. It was generally believed that the members of a Cooper pair are bound together because of interactions between the electrons and the lattice vibrations (phonons), but it is un-

likely that these interactions are strong enough to explain transition temperatures as high as 90 K. Most experts believe that interactions among the electrons generate high-temperature superconductivity. The details of this interaction are difficult to treat theoretically because the motions of the electrons are strongly correlated with each other and because magnetic phenomena play an important part in determining the microscopic properties of these materials. These strong correlations and magnetic properties may be responsible for unusual temperature dependencies of the electric resistivity ρ and Hall coefficient R_H in the normal state (*i.e.*, above T_c). (For a discussion of the Hall effect, see ELECTRICITY AND MAGNETISM: *Magnetism: Magnetic forces*.) It is observed that at temperatures above T_c the electric resistivity, although higher for superconductors than for typical metals in the normal state, is roughly proportional to the temperature T , an unusually weak temperature dependence. Measurements of R_H show it to be significantly temperature-dependent in the normal state (sometimes proportional to $1/T$) rather than being roughly independent of T , which is the case for ordinary materials.

Films of the new materials can carry currents in the superconducting state that are large enough to be of importance in making many devices. Possible applications of the high-temperature superconductors in thin-film or bulk form include the construction of computer parts (logic devices, memory elements, switches, and interconnects), oscillators, amplifiers, particle accelerators, highly sensitive devices for measuring magnetic fields, voltages, or currents, magnets for medical magnetic-imaging devices, magnetic energy-storage systems, levitated passenger trains for high-speed travel, motors, generators, transformers, and transmission lines. The principal advantages of these superconducting devices would be their low power dissipation, high operating speed, and extreme sensitivity.

Equipment made with the high-temperature superconductors would also be more economical to operate because such materials can be cooled with inexpensive liquid nitrogen (boiling point, 77 K) rather than with costly liquid helium (boiling point, 4.2 K). The ceramics have problems, however, which must be overcome before useful devices can be made from them. These problems include brittleness, instabilities of the materials in some chemical environments, and a tendency for impurities to segregate at surfaces and grain boundaries, where they interfere with the flow of high currents in the superconducting state. (D.M.Gi.)

Superfluidity

The phenomenon of superfluidity has so far been directly observed only in the liquid forms of the two stable isotopes of helium, ^4He and ^3He , both of which remain liquid under their own vapour pressure down to the lowest temperatures reached thus far. In the case of the more abundant isotope, ^4He , superfluidity occurs at temperatures below the so-called lambda transition, which is 2.17 K. It is so named because the graph of the specific heat versus the temperature near this point has a shape resembling the Greek letter lambda. For ^3He , superfluidity is observed only at temperatures far closer to absolute zero, below 3×10^{-3} K. It is believed that superfluidity may occur in some other systems, such as neutron stars, but the evidence in these cases is much less direct.

The most spectacular signature of the transition of liquid ^4He into the superfluid phase is the sudden onset of the ability to flow without apparent friction through capillaries so small that any ordinary liquid (including ^4He itself above the lambda transition) would be clamped by its viscosity; thus, a vessel that was "helium-tight" in the so-called normal phase (*i.e.*, above the lambda temperature) might suddenly spring leaks below it. Related phenomena observed in the superfluid phase include the ability to sustain persistent currents in a ring-shaped container; the phenomenon of film creep, in which the liquid flows without apparent friction up and over the side of a bucket containing it; and a thermal conductivity that is millions of times its value in the normal phase and greater than that of the best metallic conductors. Another property is

Applica-
tions
of high-
tempera-
ture super-
conductors

Superfluid
phenom-
ena

less spectacular but is extremely significant for an understanding of the superfluid phase: if the liquid is cooled through the lambda transition in a bucket that is slowly rotating, then, as the temperature decreases toward absolute zero, the liquid appears gradually to come to rest with respect to the laboratory even though the bucket continues to rotate. This nonrotation effect is completely reversible; the apparent velocity of rotation depends only on the temperature and not on the history of the system. Most of these phenomena also have been observed in the superfluid phase of liquid ^3He , though in somewhat less spectacular form.

It is thought that there is a close connection between the phenomena of superfluidity and superconductivity; indeed, from a phenomenological point of view superconductivity is simply superfluidity occurring in an electrically charged system. Thus, the frictionless flow of superfluid ^4He through narrow capillaries parallels the frictionless carrying of electric current by the electrons in a superconductor, and the ability of helium to sustain circulating mass currents in a ring-shaped container is closely analogous to the persistence of electric currents in a superconducting ring. Less obviously, it turns out that the nonrotation effect is the exact analogue of the Meissner effect in superconductors (see above *Superconductivity*). Many other characteristic features of superconductivity, such as the existence of vortices and the Josephson effect, have been observed in the superfluid phases of both ^4He and ^3He .

The accepted theoretical understanding of superfluidity (or superconductivity) is based on the idea that an extremely large number of atoms (or electrons) show identical, and moreover essentially quantum mechanical, behaviour; that is to say, the system is described by a single, coherent, quantum mechanical wave function. A single electron in an atom cannot rotate around the nucleus in any arbitrary orbit; rather, quantum mechanics requires that it rotate in such a way that its angular momentum is quantized so as to be a multiple (including zero) of $h/2\pi$, where h is Planck's constant. This is the origin of, for example, the phenomenon of atomic diamagnetism. Similarly, a single atom (or molecule) placed in a ring-shaped container is allowed by quantum mechanics to travel around the ring with only certain definite velocities, including zero. In an ordinary liquid such as water, the thermal disorder ensures that the atoms (or molecules) are distributed over the different (quantized) states available to them in such a way that the average velocity is not quantized; thus, when the container rotates and the liquid is given sufficient time to come into equilibrium, it rotates along with the container in accordance with everyday experience.

In a superfluid system the situation is quite different. The simpler case is that of ^4He , a liquid consisting of atoms that have total spin angular momentum equal to zero and whose distribution between their possible states is therefore believed to be governed by a principle known as Bose statistics. A gas of such atoms without interactions between them would undergo, at some temperature T_0 , a phenomenon known as Bose condensation; below T_0 a finite fraction of all the atoms occupy a single state, normally that of lowest energy, and this fraction increases toward one as the temperature falls toward absolute zero. These atoms are said to be condensed. It is widely believed

that a similar phenomenon should also occur for a liquid such as ^4He , in which the interaction between atoms is quite important, and that the lambda transition of ^4He is just the onset of Bose condensation. (The reason that this phenomenon is not seen in other systems of spin-zero atoms such as neon-22 is simply that, as the temperature is lowered, freezing occurs first.) If this is so, then, for temperatures below the lambda transition, a finite fraction of all the atoms must decide cooperatively which one of the possible quantized states they will all occupy. In particular, if the container is rotating at a sufficiently slow speed, these condensed atoms will occupy the non-rotating state—*i.e.*, they will be at rest with respect to the laboratory—while the rest will behave normally and will distribute themselves in such a way that on average they rotate with the container. As a result, as the temperature is lowered and the fraction of condensed atoms increases, the liquid will appear gradually to come to rest with respect to the laboratory (or, more accurately, to the fixed stars). Similarly, when the liquid is flowing through a small capillary, the condensed atoms cannot be scattered by the walls one at a time since they are forced by Bose statistics to occupy the same state. They must be scattered, if at all, simultaneously. Since this process is extremely improbable, the liquid, or more precisely the condensed fraction of it, flows without any apparent friction. The other characteristic manifestations of superfluidity can be explained along similar lines.

The idea of Bose condensation is not directly applicable to liquid ^3He , because ^3He atoms have spin angular momentum equal to $1/2$ (in units of $h/2\pi$) and their distribution among states is therefore believed to be governed by a different principle, known as Fermi statistics. It is believed, however, that in the superfluid phase of ^3He the atoms, like the electrons in a superconductor (see above *Superconductivity*), pair off to form Cooper pairs—a sort of quasimolecular complex—which have integral spin and therefore effectively obey Bose rather than Fermi statistics. In particular, as soon as the Cooper pairs are formed, they undergo a sort of Bose condensation, and subsequently the arguments given above for ^4He apply equally to them. As in the case of the electrons in superconductors, a finite energy, the so-called energy gap Δ , is necessary to break up the pairs (or at least most of them), and as a result the thermodynamics of superfluid ^3He is quite similar to that of superconductors. There is one important difference between the two cases. Whereas in a classic superconductor the electrons pair off with opposite spins and zero total angular momentum, making the internal structure of the Cooper pairs rather featureless, in ^3He the atoms pair with parallel spins and nonzero total angular momentum, so that the internal structure of the pairs is much richer and more interesting. One manifestation of this is that there are three superfluid phases of liquid ^3He , called *A*, *B*, and *A*₁, which are distinguished by the different internal structures of the Cooper pairs. The *B* phase is in most respects similar to a classic superconductor, whereas the *A* (and *A*₁) phase is strongly anisotropic in its properties and has an energy gap that actually vanishes for some directions of motion. As a result, some of the superfluid properties of the *A* and *A*₁ phases are markedly different from those of ^4He or $^3\text{He-B}$ and are indeed unique among known physical systems. (A.J.L.)

Superfluid
phases
of ^3He

Bose
conden-
sation

HIGH-PRESSURE PHENOMENA

Matter undergoes significant changes in physical, chemical, and structural characteristics when subjected to high pressure. Pressure thus serves as a versatile tool in materials research, and it is especially important in the investigation of the rocks and minerals that form the deep interior of the Earth and other planets. Pressure, defined as a force applied to an area, is a thermochemical variable that induces physical and chemical changes comparable to the more familiar effects of temperature. Liquid water, for example, transforms to solid ice when cooled to temperatures below 0°C (32°F), but ice can also be

produced at room temperature by compressing water to pressures roughly 10,000 times above atmospheric pressure. Similarly, water converts to its gaseous form at high temperature or at low pressure.

In spite of the superficial similarity between temperature and pressure, these two variables are fundamentally different in the ways they affect a material's internal energy. Temperature variations reflect changes in the kinetic energy and thus in the entropy of vibrating atoms. Increased pressure, on the other hand, alters the electron interaction energy of atomic bonds by forcing atoms closer together in

a smaller volume. Pressure thus serves as a powerful probe of atomic interactions and chemical bonding. Furthermore, pressure is an important tool for synthesizing dense structures, including superhard materials, novel solidified gases and liquids, and mineral-like phases suspected to occur deep within the Earth and other planets.

Units of pressure

Numerous units for measuring pressure have been introduced and, at times, are confused in the literature. The atmosphere (atm: approximately 1.034 kilograms per square centimetre [14.7 pounds per square inch], equivalent to the weight of about 760 millimetres [30 inches] of mercury) and the bar (equivalent to one kilogram per square centimetre) are often cited. Coincidentally, these units are almost identical (1 bar = 0.987 atm). The pascal, defined as one newton per square metre (1 Pa = 0.00001 bar), is the official SI (Système International d'Unités) unit of pressure. Nevertheless, the pascal has not gained universal acceptance among high-pressure researchers, perhaps because of the awkward necessity of using the gigapascal (1 GPa = 10,000 bars) and terapascal (1 TPa = 10,000,000 bars) in describing high-pressure results.

In everyday experience, greater-than-ambient pressures are encountered in, for example, pressure cookers (about 1.5 atm), pneumatic automobile and truck tires (usually 2 to 3 atm), and steam systems (up to 20 atm). In the context of materials research, however, "high pressure" usually refers to pressures in the range of thousands to millions of atmospheres.

Studies of matter under high pressure are especially important in a planetary context. Objects in the deepest trench of the Pacific Ocean are subjected to about 0.1 GPa (roughly 1,000 atm), equivalent to the pressure beneath a three-kilometre column of rock. The pressure at the centre of the Earth exceeds 300 GPa, and pressures inside the largest planets—Saturn and Jupiter—are estimated to be roughly 2 and 10 TPa, respectively. At the upper extreme, pressures inside stars may exceed 1,000,000,000 TPa.

Producing high pressure

Scientists study materials at high pressure by confining samples in specially designed machines that apply a force to the sample area. Prior to 1900 these studies were conducted in rather crude iron or steel cylinders, usually with relatively inefficient screw seals. Maximum laboratory pressures were limited to about 0.3 GPa, and explosions of the cylinders were a common and sometimes injurious occurrence. Dramatic improvements in high-pressure apparatuses and measuring techniques were introduced by the American physicist Percy Williams Bridgman of Harvard University in Cambridge, Mass. In 1905 Bridgman discovered a method of packing pressurized samples, including gases and liquids, in such a way that the sealing gasket always experienced a higher pressure than the sample under study, thereby confining the sample and reducing the risk of experimental failure. Bridgman not only routinely attained pressures above 30,000 atm, but he also was able to study fluids and other difficult samples.

LARGE-VOLUME APPARATUSES

Sustained high pressures and temperatures are now commonly produced in massive presses that focus large forces (up to thousands of tons) through two or more strong anvils to compress a sample. The simplest of these devices, introduced by Bridgman in the 1930s, employs two tapered anvils that squeeze the sample like a vise (see Figure 52). Although capable of very high pressures—in excess of 50 GPa in designs with sufficient lateral anvil support—the axial force of the squeezer tends to deform samples into extremely flattened, highly strained disks.

The piston-in-cylinder design, in use for more than a century, incorporates a strong metal or carbide piston that is rammed into a sample-confining cylinder. In principle, the piston can be quite long, so a piston-cylinder design can accommodate a much larger volume of sample than the squeezer, depending on the dimensions of the sample-holding cylinder. These devices are rarely used at pressures above about 10 GPa owing to the likelihood of lateral failure (namely, explosive bursting) of the metal cylinder.

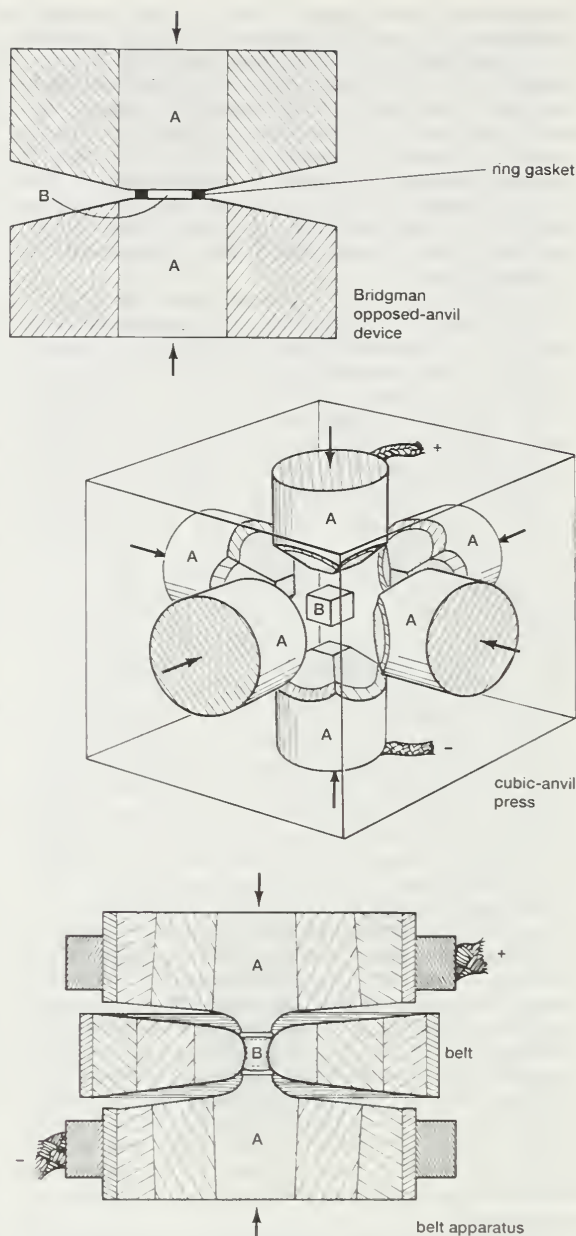


Figure 52: High-pressure apparatuses. In each device, (A) anvils of carbide (stippled) and steel compress (B) a sample. Electric leads (+ and -) provide heating capability.

Adapted from A.A. Giardini and J.E. Tydings, "Diamond Synthesis: Observations on the Mechanism of Formation," *The American Mineralogist*, vol. 47, Nov.-Dec. 1962, pp. 1393-1412, fig. 1 (a & e), copyright by the Mineralogical Society of America.

The belt apparatus, invented in 1954 by the scientist Tracy Hall of the General Electric Company for use in the company's diamond-making program, incorporates features of both opposed-anvil and piston-cylinder designs (see Figure 52). Two highly tapered pistonlike anvils compress a sample that is confined in a torus, much like a cylinder open at both ends. Hundreds of belt-type devices are in use worldwide in diamond synthesis.

Many high-pressure researchers now employ split-sphere or multianvil devices, which compress a sample uniformly from all sides. Versions with six anvils that press against the six faces of a cube-shaped sample (see Figure 52) or with eight anvils that compress an octahedral sample are in widespread use. Unlike the simple squeezer, piston-cylinder, and belt apparatuses, multianvil devices can compress a sample uniformly from all sides, while achieving a pressure range with an upper limit of at least 30 GPa. All these types of high-pressure apparatuses can be fitted with a resistance heater, typically a sample-surrounding cylinder of graphite or another electrically conducting heating element, for studies at temperatures up to 2,000° C.

Simple anvil devices

THE DIAMOND-ANVIL CELL

The diamond-anvil pressure cell, in which two gem-quality diamonds apply a force to the sample, revolutionized high-pressure research (see Figure 53). The diamond-anvil cell was invented in 1958 almost simultaneously by workers at the National Bureau of Standards in Washington, D.C., and at the University of Chicago. The diamond-cell design represented a logical outgrowth of Bridgman's simple squeezer, but it had one significant advantage over all other high-pressure apparatuses. Diamond, while extremely strong, is also transparent to many kinds of electromagnetic radiation, including gamma rays, X rays, visible light, and much of the infrared and ultraviolet region. The diamond cell thus provided the first opportunity for high-pressure researchers to observe visually the effects of pressure, and it allowed convenient access for many kinds of experimental techniques, notably X-ray diffraction, Mössbauer (gamma-ray), infrared, and Raman spectroscopies, and other optical spectroscopies.

The utility of the diamond cell was greatly enhanced when Alvin Van Valkenburg, one of the original diamond-cell inventors at the National Bureau of Standards, placed a thin metal foil gasket between the two diamond-anvil faces. Liquids and other fluid samples could thus be confined in a sample chamber defined by the cylindrical gasket wall and flat diamond ends. In 1963 Van Valkenburg became the first person to observe water, alcohol, and other liquids crystallize at high pressure. The gasketed geometry also permitted for the first time X-ray and optical studies of uncrushed single crystals that were hydrostatically pressurized by a fluid medium.

The diamond-anvil cell holds all records for sustained high pressures. The 100 GPa (megabar) mark was surpassed in December 1975 by the geophysicists Ho-kwang Mao and Peter M. Bell, both of the Geophysical Laboratory of the Carnegie Institution of Washington, in Washington, D.C., where they subsequently attained diamond-cell pressures of approximately 300 GPa. Heating of diamond-cell samples, with both resistance heaters and lasers, has extended accessible pressure-temperature conditions to those that prevail in most of the solid Earth.

The highest transient laboratory pressures are generated with high-velocity projectiles that induce extreme shock pressures (which often reach many millions of atmospheres) for times on the order of one microsecond. Shock waves generated by explosions or gas-propelled projectiles induce dramatic changes in physical properties, as well as rapid polymorphic transformations. Carefully timed intense pulses of X rays or laser light can be used to probe these transient environments. While dynamic high-pressure studies are limited by the difficulty of making precise measurements in such short time periods, these shock techniques have provided insights into changes in atomic structure and properties that occur at extreme conditions. Explosive shock compression has also become an important tool for the synthesis of microcrystalline diamond, which is employed in the polishing of gemstones and other hard materials.

Physical and chemical effects of high pressure

The principal effect of high pressure, observed in all materials, is a reduction in volume and a corresponding shortening of mean interatomic distances. Coincident with these structural modifications are numerous changes, often dramatic, in physical properties.

In four decades of high-pressure research, Bridgman, whose work was honoured by the 1946 Nobel Prize for Physics, documented effects of pressure on electric conductivity, thermal conductivity, viscosity, melting, reaction kinetics, and other material properties. Pressure was found to induce both continuous and discontinuous changes in matter. Bridgman and others observed smoothly varying trends in properties such as electric conductivity or volume versus pressure for most materials. Some substances, however, displayed sharp, reproducible discontinuities in these properties at specific pressures. Dramatic sudden drops in the electric resistance and volume of bismuth, lead, and other metals were carefully documented and provided

Advantages
of
diamond
cells

Contributions
of
Bridgman

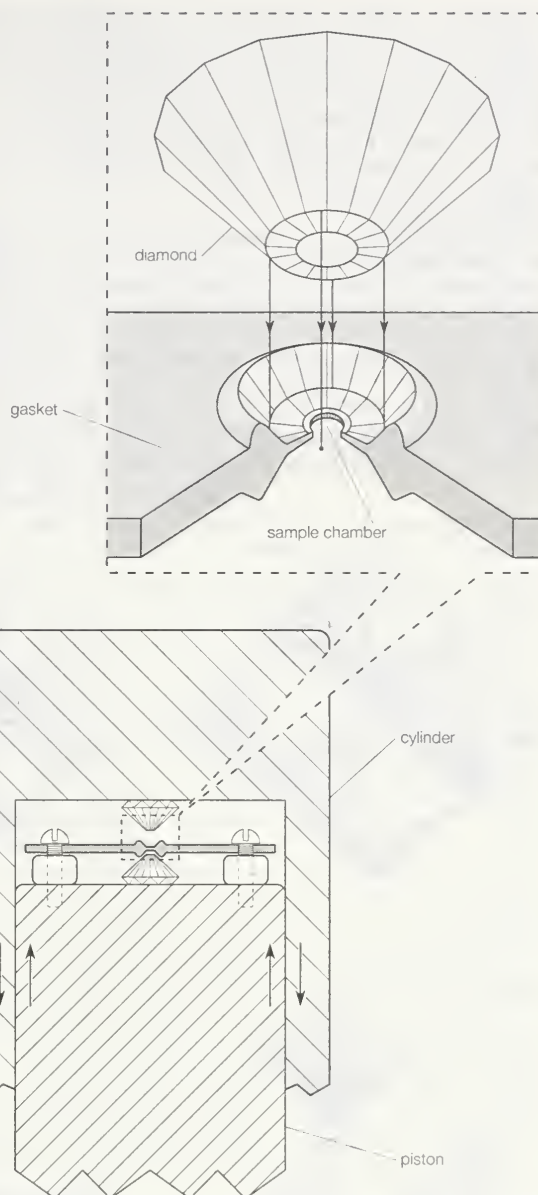


Figure 53: The diamond-anvil cell and gasket.

From H.K. Mao and P.M. Bell, "Design of the Diamond-Window, High-Pressure Apparatus for Cryogenic Experiments," *Annual Report of the Director, Geophysical Laboratory, 1978-79*, reprinted from *Papers from the Geophysical Laboratory, Carnegie Institution of Washington Year Book 78*

Bridgman with a useful internal pressure standard for his experiments. These experiments also demonstrated the effectiveness of pressure for studying continuous changes in properties (under uniform compression) and discontinuous changes (phase transitions).

PHASE TRANSITIONS

Under sufficiently high pressure, every material is expected to undergo structural transformations to denser, more closely packed atomic arrangements. At room temperature, for example, all gases solidify at pressures not greater than about 15 GPa. Molecular solids like water ice (H_2O) and carbon tetrachloride (CCl_4) often undergo a series of structural transitions, characterized by successively denser arrangements of molecular units.

A different transition mode is observed in oxides, silicates, and other types of ionic compounds that comprise most rock-forming minerals. In these materials, metal or semimetal atoms such as magnesium (Mg) or silicon (Si) are surrounded by regular tetrahedral or octahedral arrangements of four or six oxygen (O) atoms, respectively. High-pressure phase transitions of such minerals often involve a structural rearrangement that increases the number of oxygen atoms around each central cation. The common mineral quartz (SiO_2), for example, contains four-coordinated silicon at low pressure, but it transforms

to the dense stishovite form with six-coordinated silicon at about 8 GPa. Similarly, the pyroxene mineral with formula $MgSiO_3$ at room pressure contains magnesium and silicon in six- and four-coordination, respectively, but the pyroxene transforms to the perovskite structure with eight-coordinated magnesium and six-coordinated silicon above 25 GPa. Each of these high-pressure phase transitions results in a denser structure with increased packing efficiency of atoms.

High-pressure metallization

The British scientist J.D. Bernal predicted in 1928 that all matter should ultimately become metallic at sufficient pressure, as the forced overlap of electron orbitals induces electron delocalization. High-pressure transformations from insulator to metal were first observed in iodine, silicon, germanium, and other elements by the American chemist Harry G. Drickamer and his coworkers at the University of Illinois at Urbana-Champaign in the early 1960s. Subsequently, metallization has been documented in several more elements (including the gases xenon and

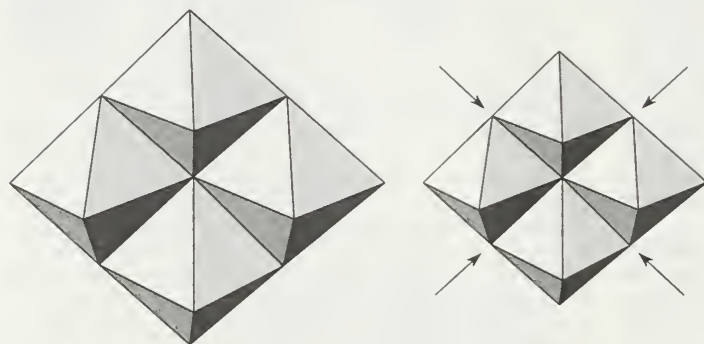
oxygen), as well as in numerous molecular, ionic, and covalent chemical compounds. The effort to metallize the element hydrogen at a predicted pressure of several million atmospheres remains a significant challenge in experimental physics.

COMPRESSION

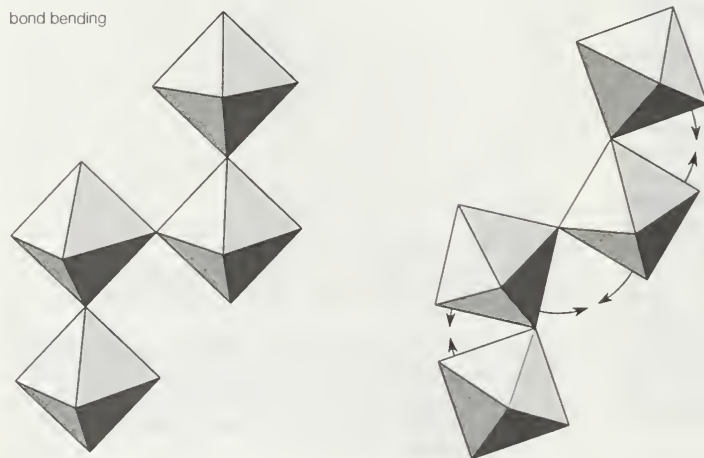
High-pressure X-ray crystallographic studies of atomic structure reveal three principal compression mechanisms in solids: bond compression, bond-angle bending, and intermolecular compression; they are illustrated in Figure 54. Bond compression—*i.e.*, the shortening of interatomic distances—occurs to some extent in all compounds at high pressure. The magnitude of this effect has been shown both theoretically and empirically to be related to bond strength. Strong covalent carbon-carbon bonds in diamond experience the lowest percentage of compression: roughly 0.07 percent per GPa. Similarly, ionic bonds between highly charged cations and anions, such as bonds between

From R.M. Hazen and L.W. Finger, "Crystals at High Pressure," copyright © 1985 by Scientific American, Inc., all rights reserved

bond compression



bond bending



intermolecular compression

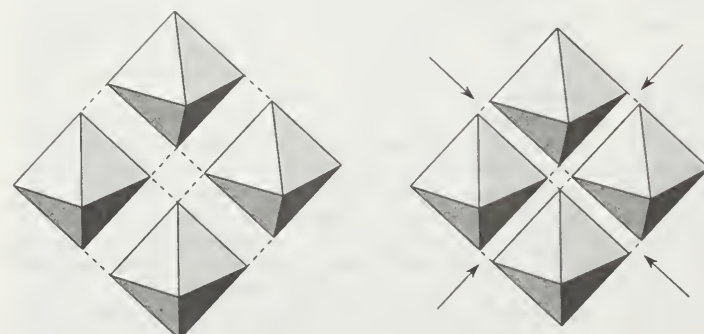


Figure 54: Three compression mechanisms in crystals.

Si^{4+} and O^{2-} in silicates, are relatively incompressible (less than 0.2 percent per GPa). Relatively weak bonds in alkali halides, on the other hand, display bond compressibilities that often exceed 5.0 percent per GPa.

Many common materials display different bonding characteristics in different directions; this occurs notably in layered compounds (*e.g.*, graphite and layered silicates such as micas) and in chain compounds (*e.g.*, many polymeric compounds and chain silicates, including some varieties of asbestos). The strong dependence of bond compression on bond strength thus commonly leads to anisotropies—that is, significant differences in compression in different crystal directions. In many layered-structure silicates, such as mica, in which relatively strong and rigid layers containing magnesium-oxygen, aluminum-oxygen, and silicon-oxygen bonds alternate with weaker layers containing alkali cations, compressibility is five times greater perpendicular to the layers than within the layers. This differential compressibility and the associated stresses that develop in a high-pressure geologic environment contribute to the development of dramatic layered textures in mica-rich rocks such as schist.

Many common ionic compounds, including the rock-forming minerals quartz, feldspar, garnet, zeolite, and perovskite (the high-pressure MgSiO_3 form of which is thought to be the Earth's most abundant mineral), are composed of corner-linked clusters—or frameworks—of atomic polyhedrons. A polyhedron consists of a central cation, typically silicon or aluminum in common minerals, surrounded by a regular tetrahedron or octahedron of four or six oxygen atoms, respectively. In framework structures every oxygen atom is bonded to two tetrahedral or octahedral cations, resulting in a three-dimensional polyhedral network. In these materials significant compression can occur by bending the metal-oxygen-metal bond angles between the polyhedrons. The volume change resulting from this bending, and the associated collapse of interpolyhedral spaces, is typically an order of magnitude greater than compression due to bond-length changes alone. Framework structures, consequently, are often much more compressible than structures with only edge- or face-sharing polyhedrons, whose compression is attributable predominantly to bond shortening.

Molecular solids—including ice, solidified gases such as solid oxygen (O_2), hydrogen (H_2), and methane (CH_4), and virtually all organic compounds—consist of an array of discrete, rigid molecules that are linked to one another by weak hydrogen bonds and van der Waals forces. Compression in these materials generally occurs by large decreases in intermolecular distances (often approaching 10 percent per GPa), in contrast to minimal intramolecular compression. Differences in the intermolecular versus intramolecular compression mechanisms lead in some cases to significantly anisotropic compression. Graphite, the low-pressure layered form of elemental carbon in which the “molecules” are continuous two-dimensional sheets, exhibits perhaps the most extreme example of this phenomenon. Carbon-carbon bonds within graphite layers compress only 0.07 percent per GPa (similar to C-C bond compression in diamond), while interlayer compression, dominated by van der Waals forces acting between carbon sheets, is approximately 45 times greater.

EFFECTS ON ELECTRIC AND MAGNETIC PROPERTIES

The measurement of electric and magnetic properties of materials in a high-pressure environment entails considerable experimental difficulties, especially those associated with attaching leads to pressurized samples or detecting small signals from the experiment. Nevertheless, electric conductivities of numerous materials at high pressures have been documented. The principal classes of solids—insulators, semiconductors, metals, and superconductors—are distinguished on the basis of electric conductivity and its variation with temperature (see above *Solid state*). Insulators, which include most rock-forming oxides and silicates, have been investigated extensively by geophysicists concerned primarily with the behaviour and properties of deep-earth rocks and minerals at extreme conditions. Indeed, it was once hoped that laboratory constraints on

such properties could be tied to known values of the Earth's electric and magnetic properties and thus constrain the composition and temperature gradients of Earth models. It appears, however, that small variations in mineral composition (*e.g.*, the ratio of ferrous to ferric iron) as well as defect properties can play a role orders of magnitude greater than that of pressure alone.

Properties of semiconductors are highly sensitive to pressure, because small changes in structure can result in large changes in electronic properties. The metallizations of silicon and germanium, which are accompanied by an orders-of-magnitude increase in electric conductivity, represent extreme cases of such changes. While simple metals display a general trend of increased conductivity with increased pressure, there are many exceptions. Calcium and strontium exhibit maxima in electric conductivity at 30 and 4 GPa, respectively, while barium and arsenic display both maxima and minima with increasing pressure. Ionic conductors, on the other hand, generally experience decreased electric conductivity at high pressure owing to the collapse of ion pathways.

Pressure has been found to be a sensitive probe of the effects of structure on superconductivity, because the structural changes brought about by pressure often have a significant effect on the critical temperature. In simple metals, pressure tends to decrease the critical temperature, eventually suppressing superconductivity altogether. In some organic superconductors, on the other hand, superconductivity appears only at high pressure (and temperatures near absolute zero). In several of the layered copper-oxide high-temperature superconductors, pressure has a strong positive effect on critical temperature; this phenomenon led to the synthesis of new varieties of superconductors in which smaller cations are used to mimic the structural effect of pressure.

The first measurements of magnetic properties at high pressure were conducted on samples in a diamond-anvil cell using Mössbauer spectroscopy, which is a technique that can probe the coupling of a magnetic field with the nuclear magnetic dipole. High-pressure ferromagnetic-to-paramagnetic transitions were documented in iron metal and in magnetite (Fe_3O_4), while Curie temperatures (*i.e.*, the temperature above which the ferromagnetic properties of a material cease to exist) in several metallic elements were found to shift slightly. Subsequent research has employed high-pressure devices constructed of nonmagnetic beryllium-copper alloys, which were developed for research on samples subjected to strong magnetic fields.

Applications

DIAMOND MAKING

While modest pressures (less than 1,000 atm) have long been used in the manufacture of plastics, in the synthesis of chromium dioxide for magnetic recording tape, and in the growth of large, high-quality quartz crystals, the principal application of high-pressure materials technology lies in the synthesis of diamond and other superhard abrasives. Approximately 100 tons of synthetic diamond are produced each year—a weight comparable to the total amount of diamond mined since biblical times. For centuries diamonds had been identified only as an unusual mineral found in river gravels; scientists had no clear idea about their mode of origin until the late 1860s, when South African miners found diamond embedded in its native matrix, the high-pressure volcanic rock called kimberlite. Efforts to make diamond by subjecting graphitic carbon to high pressure began shortly after that historic discovery.

Prior to the work of Bridgman, sufficient laboratory pressures for driving the graphite-to-diamond transition had not been achieved. Bridgman's opposed-anvil device demonstrated that the necessary pressures could be sustained, but high temperatures were required to overcome the kinetic barrier to the transformation. Following World War II, several industrial laboratories, including Allmanna Svenska Elektriska Aktiebolaget (ASEA) in Sweden and Norton Company and General Electric in the United States, undertook major efforts to develop a commercial process. Diamond was first synthesized in a reproducible,

Effects on superconductors

Anisotropic compression of graphite

commercially viable experiment in December 1954, when Tracy Hall, working for General Electric, subjected a mixture of iron sulfide and carbon to approximately 6 GPa and 1,500° C in a belt-type apparatus. General Electric employees soon standardized the processes and discovered that a melted ferrous metal, which acts as a catalyst, is essential for diamond growth at these conditions.

EARTH SCIENCE

Diamond-making techniques have been embraced by Earth scientists in their efforts to simulate conditions in the Earth's deep interior. Of special significance were the high-pressure syntheses of two new forms of silicates. In 1960 Sergei Stishov, while at the Institute of High-Pressure Physics in Moscow, subjected ordinary beach sand (composed of the mineral quartz SiO₂) to more than 8 GPa of pressure and high temperatures. The form of silica that he produced was approximately 62 percent denser than quartz and was the first known high-pressure compound to contain silicon in six-coordination rather than the four-coordination found in virtually all crustal minerals. The natural occurrence of this new synthetic material was confirmed within a few weeks by careful examination of shocked material from Meteor Crater, Ariz., U.S. The mineral was named stishovite.

In 1974 a second high-pressure discovery revolutionized geologists' understanding of deep-earth mineralogy when Lin-gun Liu of the Australian National University used a diamond-anvil cell to synthesize silicate perovskite, a dense form of the common mineral enstatite, MgSiO₃. Subsequent studies by Liu revealed that many of the minerals believed to constitute the deep interior of the Earth transform to the perovskite structure at lower mantle conditions—an observation that led him to propose that silicate perovskite is the Earth's most abundant mineral, perhaps accounting for more than half of the planet's volume. (R.M.Ha.)

BIBLIOGRAPHY

Phase changes. DANIEL S. BARKER, *Igneous Rocks* (1983), ch. 3, "Phase Relations," pp. 24–57, presents a clear summary of the interpretation of petrologic phase diagrams. ERNEST G. EHLERS, *The Interpretation of Geological Phase Diagrams* (1972, reprinted 1987), provides step-by-step nonmathematical procedures for understanding both simple and complex phase diagrams. ERNEST G. EHLERS and HARVEY BLATT, *Petrology: Igneous, Sedimentary, and Metamorphic* (1982), ch. 2, "Experiments with Molten Silicates: Unary and Binary Systems," pp. 43–73, provides an introduction to phase equilibria of petrologic systems. DONALD W. HYNDMAN, *Petrology of Igneous and Metamorphic Rocks* (1972), in the second half of ch. 1, "Environment and Materials," pp. 15–30, summarizes the interpretation of some important petrologic phase diagrams. W.G. ERNST, *Petrologic Phase Equilibria* (1976), is a concise introduction to phase equilibria; some knowledge of thermodynamics on the part of the reader would be helpful. (E.G.E.)

Solid state. General works include LAWRENCE H. VAN VLACK, *Elements of Materials Science and Engineering*, 6th ed. (1989), an elementary textbook; CHARLES A. WERT and ROBB M. THOMSON, *Physics of Solids*, 2nd ed. (1970), an intermediate-level text; CHARLES KITTEL, *Introduction to Solid State Physics*, 6th ed. (1986), the standard college textbook; NEIL W. ASHCROFT and N. DAVID MERMIN, *Solid State Physics* (1976), an advanced textbook; GEORGE E. BACON, *The Architecture of Solids* (1981), an introduction to bonding and structure; and LINUS PAULING, *The Nature of the Chemical Bond and the Structure of Molecules and Crystals*, 3rd ed. (1960, reissued 1989), the classic reference work on chemical bonding.

Crystalline solids: A readable treatment of crystal structure is C.S. BARRETT and T.B. MASSALSKI, *Structure of Metals: Crystallographic Methods, Principles, and Data*, 3rd rev. ed. (1980). RALPH W.G. WYCKOFF, *Crystal Structures*, 2nd ed., 6 vol. (1963–71), compiles information on known structures. Phase diagrams of binary alloys are compiled in T.B. MASSALSKI (ed.), *Binary Alloy Phase Diagrams*, 2nd ed., 3 vol. (1990); and in WILLIAM G. MOFFATT (ed.), *The Handbook of Binary Phase Diagrams*, 4 vol. (1976–), published in a regularly updated looseleaf format. Traditional crystal growing is discussed in J.C. BRICE, *The Growth of Crystals from the Melt* (1965); ALAN COTTRELL, *An Introduction to Metallurgy*, 2nd ed. (1975); and MERTON C. FLEMINGS, *Solidification Processing* (1974). Epitaxy is discussed by GERALD B. STRINGFELLOW, *Organometallic Vapor Phase Epitaxy: Theory and Practice* (1989); and SEYMOUR P. KELLER (ed.), *Materials, Properties, and Preparation* (1980). WILLIAM SHOCK-

LEY, *Electrons and Holes in Semiconductors* (1950, reprinted 1976), is an introduction; while RICHARD DALVEN, *Introduction to Applied Solid State Physics*, 2nd ed. (1990), is an intermediate-level text on semiconductors and semiconductor devices. A wealth of information is available in H.H. LANDOLT and R. BÖRNSTEIN, *Zahlenwerte und Funktionen aus Naturwissenschaften und Technik*, new series, ed. by K.H. HELLEWEGE, known as *Landolt-Börnstein*, with a parallel title in English, *Numerical Data and Functional Relationships in Science and Technology*, reflecting the two languages of the work; group III, *Kristall- und Festkörperphysik* (or *Crystal and Solid State Physics*), contains compiled measurements of semiconductor properties in vol. 17a–g (1982–85) and resistivity data on all metals and alloys in vol. 15a, pp. 1–289 (1982), and vol. 15b, pp. 1–47 (1985). Books on magnetism are D.H. MARTIN, *Magnetism in Solids* (1967); JOHN E. THOMPSON, *The Magnetic Properties of Materials* (1968); and KENNETH HOPE STEWART, *Ferromagnetic Domains* (1954). Fullerenes are reviewed in a special issue of *Accounts of Chemical Research*, vol. 25, no. 3 (March 1992). (G.D.Ma.)

Quasicrystals: Introductions to the topic are available in DAVID R. NELSON, "Quasicrystals," *Scientific American*, 255(2):43–51 (August 1986); PETER W. STEPHENS and ALAN I. GOLDMAN, "The Structure of Quasicrystals," *Scientific American*, 264(4):44–47, 50–53 (April 1991); and P.J. STEINHARDT, "Icosahedral Solids: A New Phase of Matter?," *Science*, 238(4831):1242–47 (Nov. 27, 1987). MARTIN GARDNER, "Mathematical Games," *Scientific American*, 236(1):110–112, 115–121 (January 1977), discusses Penrose tilings and their remarkable properties. More technically detailed works are D.P. DIVENCENZO and P.J. STEINHARDT (eds.), *Quasicrystals: The State of the Art* (1991); and the series *Aperiodicity and Order*, ed. by MARKO V. JARIĆ (1988–).

Liquid crystals: The history of the field is surveyed by H. KELKER, "History of Liquid Crystals," *Molecular Crystals and Liquid Crystals*, 21(1 and 2):1–48 (May 1973). The Nobel Prize acceptance lecture by P.G. DE GENNES, "Soft Matter," *Reviews of Modern Physics*, 64(3):645–648 (July 1992), sets liquid crystals in a broader scientific context. Discussions of special topics in liquid crystals, frequently at a level close to this article, may be found in the periodical *Condensed Matter News* (bimonthly). More technical presentations are given in P.G. DE GENNES, *The Physics of Liquid Crystals* (1974); S. CHANDRASEKHAR, *Liquid Crystals*, 2nd ed. (1992); and P.S. PERSHAN, *Structure of Liquid Crystal Phases* (1988). Applications of liquid crystals are described in E. KANEKO, *Liquid Crystal TV Displays* (1987); and J. FUNFSCHILLING, "Liquid Crystals and Liquid Crystal Displays," *Condensed Matter News*, 1:12–16 (1991). (M.W.)

Amorphous solids: A lucid introductory text accessible to a nontechnical reader is RICHARD ZALLEN, *The Physics of Amorphous Solids* (1983), with coverage of structural models for the various classes of amorphous solids as well as percolation theory, a modern paradigm for disordered systems. A classic advanced work is N.F. MOTT and E.A. DAVIS, *Electronic Processes in Non-crystalline Materials*, 2nd ed. (1979), which features many of the theoretical contributions of Nobel Laureate coauthor Mott. A text providing a thorough treatment of oxide glasses is J. ZARZYCKI, *Glasses and the Vitreous State* (1991; originally published in French, 1982). A reference work with wide coverage of recent research topics, including detailed treatment of chalcogenide glasses, is S.R. ELLIOTT, *Physics of Amorphous Materials*, 2nd ed. (1990). A comprehensive collection of detailed reviews is contained in R.W. CAHN, P. HAASSEN, and E.J. KRAMER (eds.), *Materials Science and Technology*, vol. 6, *Glasses and Amorphous Materials*, ed. by J. ZARZYCKI (1991), including coverage of glass technology, formation, and structure, oxide glasses, chalcogenide glasses, metallic glasses, polymeric glasses, and the optical, electric, and mechanical properties of glasses. Amorphous silicon is treated in detail in another work, R.A. STREET, *Hydrogenated Amorphous Silicon* (1991). (R.Z.)

Liquid state. J.N. MURRELL and E.A. BOUCHER, *Properties of Liquids and Solutions* (1982), is a short introduction to the physics and chemistry of the liquid state. ROBERT C. REID, JOHN M. PRAUSNITZ, and BRUCE E. POLING, *The Properties of Gases and Liquids*, 4th ed. (1987), focuses on the vapour-liquid transition, as opposed to the solid-liquid transition, and evaluates and illustrates techniques for estimating and correlating properties of gases and liquids, as well as tabulating the properties of 600 compounds. A. BONDI, *Physical Properties of Molecular Crystals, Liquids, and Glasses* (1968), focuses on the solid-liquid transition as opposed to the vapour-liquid transition and describes methods for the characterization of higher-molecular-weight liquids. J.S. ROWLINSON and E.L. SWINTON, *Liquids and Liquid Mixtures*, 3rd ed. (1982), provides a thorough treatment of the physics of fluids and gives some statistical mechanical theories of the equilibrium properties

of simple pure liquids and liquid mixtures; the work also contains a data bibliography and is primarily for research-oriented readers. *Dictionary of Organic Compounds*, 5th ed., 7 vol. (1982), and annual supplements, is a listing of such properties of organic compounds as chemical formula, density, and melting and boiling points; the work is useful for organic chemists in synthesizing and identifying organic compounds. GILBERT NEWTON LEWIS and MERLE RANDALL, *Thermodynamics*, 2nd ed., rev. by KENNETH S. PITZER and LEO BREWER (1961), is a classic text on the thermodynamics of pure substances and solutions, written from the point of view of chemistry. JOHN M. PRAUSNITZ, RUEDIGER N. LICHTENTHALER, and EDMUNDO GOMES DE AZEVEDO, *Molecular Thermodynamics of Fluid-Phase Equilibria*, 2nd ed. (1986), employs molecular-thermodynamic concepts useful for engineering and is written from a chemical-engineering point of view. The first three chapters of K.E. WEALE, *Chemical Reactions at High Pressures* (1967), give a description of the behaviour of pure systems and phase equilibria at very high pressures. JOHN M. PRAUSNITZ *et al.*, *Computer Calculations for Multicomponent Vapor-Liquid and Liquid-Liquid Equilibria* (1980), includes a short, annotated list of literature sources for vapour-liquid and liquid-liquid phase equilibria. (B.E.P.)

Gaseous state. The best elementary discussion of the kinetic theory of gases is T.G. COWLING, *Molecules in Motion* (1960), clearly explaining all the fundamental physical ideas without excessive mathematical manipulation. Other elementary books are JOEL H. HILDEBRAND, *An Introduction to Molecular Kinetic Theory* (1963); SIDNEY GOLDEN, *Elements of the Theory of Gases* (1964); and WALTER KAUFMANN, *Kinetic Theory of Gases* (1966). Two excellent books make greater demands on the mathematical background of the reader: JAMES JEANS, *An Introduction to the Kinetic Theory of Gases* (1940, reissued 1982), which is an abridged and slightly simplified version of the author's classic *The Dynamical Theory of Gases*, 4th ed. (1925, reissued 1954); and RICHARD D. PRESENT, *Kinetic Theory of Gases* (1958), an excellent though selective textbook. The historical literature is especially rich; the following works may be profitably consulted: J.S. ROWLINSON (ed.), *J.D. van der Waals: On the Continuity of the Gaseous and Liquid States* (1988), a translation of van der Waals's 1873 Dutch thesis with a marvelous extended introductory essay by the editor that is unique in the field and surveys modern developments in the theory of liquids and solutions; STEPHEN G. BRUSH (ed.), *Kinetic Theory*, 3 vol. (1965-72), a set of famous historical papers along with introductory commentaries and summaries by the editor; ROBERT LINDSAY (ed.), *Early Concepts of Energy in Atomic Physics* (1979), selections from famous historical papers—many of them omitted from the previous work—together with the editor's comments; and STEPHEN G. BRUSH, *The Kind of Motion We Call Heat: A History of the Kinetic Theory of Gases in the 19th Century*, 2 vol. (1976, reissued 1986), a thorough historical account without much mathematics, and *Statistical Physics and the Atomic Theory of Matter: From Boyle and Newton to Landau and Onsager* (1983), covering a much broader range and requiring a thorough scientific background. More advanced professional treatments include SYDNEY CHAPMAN and T.G. COWLING, *The Mathematical Theory of Non-uniform Gases*, 3rd ed. (1970, reprinted 1990), the acknowledged classic on the modern kinetic theory of gases based on the Enskog-Chapman approach; JOSEPH O. HIRSCHFELDER, CHARLES F. CURTISS, and R. BYRON BIRD, *Molecular Theory of Gases and Liquids* (1954, reissued with added notes 1964), a monumental compendium of detailed results on equations of state, transport properties of gases, and intermolecular forces, especially valuable as a reference; J.H. FERZIGER and H.G. KAPER, *Mathematical Theory of Transport Processes in Gases* (1972), successfully combining the best features of the previous two works; CARLO CERCIGNANI, *The Boltzmann Equation and Its Applications* (1988), by a mathematician, one of the few books that concerns itself with free-molecule gases and the transition to continuum behaviour; and FREDERICK R.W. MCCOURT *et al.*, *Nonequilibrium Phenomena in Polyatomic Gases*, 2 vol. (1990-91), an account of the extension of the Enskog-Chapman theory to include truly molecular shape effects, including the effects of external electric and magnetic fields and of surface collisions. Special topics are addressed in MARTIN KNUDSEN, *The Kinetic Theory of Gases*, 3rd ed. (1950), a series of lectures from 1933 on rarefied gas phenomena, by one of the experimental pioneers in the subject; K.E. GREW and T.L. IBBS, *Thermal Diffusion in Gases* (1952), a short monograph on one of the more intriguing special topics of the kinetic theory of gases; J.S. ROWLINSON, *The Perfect Gas* (1963), emphasizing the internal mechanics of molecules as related to the calculation of the thermodynamic properties of gases by statistical mechanics; and E.A. MASON and T.H. SPURLING, *The Virial Equation of State* (1969), emphasizing experimental measurements and the detailed connection with intermolecular forces. (E.A.M.)

Plasma state. For a general audience an early work that is still of considerable use is LEV A. ARZIMOVICH (L.A. ARTSIMOVICH), *Elementary Plasma Physics* (1965; originally published in Russian, 1963). A more recent work at the same level is YAFFA ELIEZER and SHALOM ELIEZER, *The Fourth State of Matter: An Introduction to the Physics of Plasma* (1989). A broad perspective on plasma physics is provided by NATIONAL RESEARCH COUNCIL (U.S.), PANEL ON THE PHYSICS OF PLASMAS AND FLUIDS, *Plasmas and Fluids* (1986). On the applied side, see NATIONAL RESEARCH COUNCIL (U.S.), PANEL ON PLASMA PROCESSING OF MATERIALS, *Plasma Processing of Materials: Scientific Opportunities and Technological Challenges* (1991). Descriptions of fusion-energy options with a discussion of plasma aspects are included in ROGER A. HINRICHS, *Energy* (1992); and RUTH HOWES and ANTHONY FAINBERG (eds.), *The Energy Sourcebook* (1991). A historical treatment with numerous illustrations dealing with the aurora is ROBERT H. EATHER, *Majestic Lights* (1980). More advanced mid-college-level texts include FRANCIS F. CHEN, *Introduction to Plasma Physics and Controlled Fusion*, vol. 1, *Plasma Physics* (1984); and MICHAEL C. KELLEY and RODNEY A. HEELIS (RODERICK A. HEELIS), *The Earth's Ionosphere: Plasma Physics and Electrodynamics* (1989). (M.C.K.)

Clusters. MICHAEL A. DUNCAN and DENNIS H. ROUBRAY, "Microclusters," *Scientific American*, 261(6):110-115 (December 1989), provides a general introduction and survey for nonscientists. Works presenting the results of recent research include R. STEPHEN BERRY, "When the Melting and Freezing Points Are Not the Same," *Scientific American*, 263(2):68-72, 74 (August 1990), written for nonscientists, a description of the melting and freezing of clusters and their relation to bulk melting and freezing; and several collections of conference proceedings: P. JENA, B.K. RAO, and S.N. KHANNA (eds.), *Physics and Chemistry of Small Clusters* (1987), covering a wide variety of topics within cluster science; S. SUGANO, Y. NISHINA, and S. OHNISHI (eds.), *Microclusters* (1987); G. SCOLFS (ed.), *The Chemical Physics of Atomic and Molecular Clusters* (1990), at the graduate-student level; S. SUGANO, *Microcluster Physics* (1991), accessible to scientifically literate readers; and *Zeitschrift für Physik*, part D, vol. 19 and 20 (1991) and vol. 26 (1993), the proceedings of the international conferences on small particles and inorganic clusters held in 1990 and 1992, respectively. (R.S.Be.)

Low-temperature phenomena. F.E. SIMON *et al.*, *Low Temperature Physics: Four Lectures* (1952, reprinted 1961); and K. MENDELSSOHN, *The Quest for Absolute Zero: The Meaning of Low Temperature Physics*, 2nd ed. (1977), are excellent non-technical summaries. Another account is found in ANTHONY LEGGETT, "Low Temperature Physics, Superconductivity, and Superfluidity," ch. 9 in PAUL DAVIES (ed.), *The New Physics* (1989), pp. 268-288. The *Journal of Low Temperature Physics* (semiannual) contains useful articles. Also informative are *Progress in Low Temperature Physics* (irregular), theoretical and experimental research reports and review articles, all with extensive bibliographies; K. MENDELSSOHN (ed.), *Progress in Cryogenics*, 4 vol. (1959-64), which complements the previous work, with more emphasis on applied problems and developments; FRITZ LONDON, *Superfluids*, vol. 1, *Macroscopic Theory of Superconductivity*, 2nd ed. (1961), and vol. 2, *Macroscopic Theory of Superfluid Helium* (1954, reprinted 1964), the classic presentation of theory; H.M. ROSENBERG, *Low Temperature Solid State Physics: Some Selected Topics* (1963), fairly elementary; MARSHALL SITTIG, *Cryogenics: Research and Applications* (1963); two articles from the *American Journal of Physics*, both by D.M. GINSBERG, "Resource Letter Scy-1 on Superconductivity," 32:85-89 (1964), and "Resource Letter Scy-2 on Superconductivity," 38:949-955 (1970), introductory reviews with descriptive bibliographies; P.G. DE GENNES, *Superconductivity of Metals and Alloys*, trans. from French (1966, reprinted 1989), a presentation of the foundations and applications of the theory; CHARLES G. KUPER, *An Introduction to the Theory of Superconductivity* (1968); WILLIAM E. KELLER, *Helium-3 and Helium-4* (1969), advanced-level summaries of theory; R.D. PARKS (ed.), *Superconductivity*, 2 vol. (1969); W.D. GREGORY, W.N. MATTHEWS, JR., and E.A. FELSACK (eds.), *The Science and Technology of Superconductivity*, 2 vol. (1973); MICHAEL TINKHAM, *Introduction to Superconductivity* (1975, reprinted 1980), an intermediate introduction; A.C. ROSEFINNES and F.I. RHODERICK, *Introduction to Superconductivity*, 2nd ed. (1978), basic, with an emphasis on application; GUY K. WHITE, *Experimental Techniques in Low-Temperature Physics*, 3rd ed. (1979, reissued 1987), a detailed discussion; DAVID R. TILLEY and JOHN TILLEY, *Superfluidity and Superconductivity*, 3rd ed. (1990); and DIETER VOLLHARDT and PETER WÖLFLE, *The Superfluid Phases of Helium 3* (1990).

Authoritative review articles of higher-temperature superconductors can be found in D.M. GINSBERG (ed.), *Physical Properties of High Temperature Superconductors*, 3 vol. (1989-92). Other

reports include those by M.K. WU *et al.*, "Superconductivity at 93 K in a New Mixed-Phase Y-Ba-Cu-O Compound System at Ambient Pressure," *Physical Review Letters*, 58(9):908-910 (March 2, 1987); ANIL KHURANA, "Superconductivity Seen Above the Boiling Point of Nitrogen," *Physics Today*, 40(4):17-23 (April 1987); K. ALEX MÜLLER and J. GEORG BEDNORZ, "The Discovery of a Class of High-Temperature Superconductors," *Science*, 237(4819):1133-39 (Sept. 4, 1987); and STUART A. WOLF and VLADIMIR Z. KRESIN, *Novel Superconductivity* (1987).

(A.J.L./D.M.G.)

High-pressure phenomena. P.W. BRIDGMAN, *The Physics of High Pressure*, new impression with supplement (1949), remains the major source of information on high-pressure history and technology prior to its publication. MAILA L. WALTERS, *Science and Cultural Crisis* (1990), surveys Percy Bridgman's accomplishments. ROBERT M. HAZEN, *The New Alchemists: Breaking*

Through the Barriers of High Pressure (1993), reviews the history of high-pressure research, particularly efforts to synthesize diamond. Collected papers from several biannual conferences of the International Association for High-Pressure Research are available in B. VODAR and PH. MARTEAU (eds.), *High Pressure Science and Technology*; 2 vol. (1980); C. HOMAN, R.K. MACCRONE, and E. WHALLEY (eds.), *High Pressure in Science and Technology*; 3 vol. (1984); N.J. TRAPPENIERS *et al.* (eds.), *Proceedings of the Xth AIRAPT International High Pressure Conference on Research in High Pressure Science & Technology* (1986); and W.B. HOLZAPFEL and P.G. JOHANNSEN (eds.), *High Pressure Science and Technology* (1990). Efforts in the Earth sciences are reviewed by S. AKIMOTO and M.H. MANGHNANI (eds.), *High-Pressure Research in Geophysics* (1982); and M.H. MANGHNANI and YASUHIKO SYONO (eds.), *High-Pressure Research in Mineral Physics* (1987). (R.M.Ha.)

Maxwell

James Clerk Maxwell is regarded by most modern physicists as the scientist of the 19th century who had the greatest influence on 20th-century physics; he is ranked with Sir Isaac Newton and Albert Einstein for the fundamental nature of his contributions. In 1931, at the 100th anniversary of Maxwell's birth, Einstein described the change in the conception of reality in physics that resulted from Maxwell's work as "the most profound and the most fruitful that physics has experienced since the time of Newton." The concept of electromagnetic radiation originated with Maxwell, and his field equations, based on Michael Faraday's observations of the electric and magnetic lines of force, paved the way for Einstein's special theory of relativity, which established the equivalence of mass and energy. Maxwell's ideas also ushered in the other major innovation of 20th-century physics, the quantum theory. His description of electromagnetic radiation led to the development (according to classical theory) of the ultimately unsatisfactory law of heat radiation, which prompted Max Planck's formulation of the quantum hypothesis—*i.e.*, the theory that radiant-heat energy is emitted only in finite amounts, or quanta. The interaction between electromagnetic radiation and matter, integral to Planck's hypothesis, in turn has played a central role in the development of the theory of the structure of atoms and molecules.

By courtesy of King's College London



Maxwell, engraving by G.J. Stodart.

Early life. Maxwell came from a comfortable middle-class background. The original family name was Clerk, the additional surname being added by his father after he had inherited the Middlebie estate from Maxwell ancestors. James, an only child, was born on June 13, 1831, in Edinburgh, where his father was a lawyer; his parents had married late in life, and his mother was 40 years old at his birth. Shortly afterward the family moved to Glenlair, the country house on the Middlebie estate.

His mother died in 1839 from abdominal cancer, the very disease to which Maxwell was to succumb at exactly the same age. A dull and uninspired tumb was engaged who claimed that James was slow at learning, though in fact he displayed a lively curiosity at an early age and had a phenomenal memory. Fortunately he was rescued by his aunt Jane Cay and from 1841 was sent to school at the Edinburgh Academy. Among the other pupils were his biographer Lewis Campbell and his friend Peter Guthrie Tait.

Maxwell's interests ranged far beyond the school syllabus, and he did not pay particular attention to examination

performance. His first scientific paper, published when he was only 14 years old, described a generalized series of oval curves that could be traced with pins and thread by analogy with an ellipse. This fascination with geometry and with mechanical models continued throughout his career and was of great help in his subsequent research.

At the age of 16 he entered the University of Edinburgh, where he read voraciously on all subjects and published two more scientific papers. In 1850 he went to the University of Cambridge, where his exceptional powers began to be recognized. His mathematics teacher, William Hopkins, was a well-known "wrangler maker" (a wrangler is one who takes first class honours in the mathematics examinations at Cambridge) whose students included Tait, George Gabriel (later Sir George) Stokes, William Thomson (later Lord Kelvin), Arthur Cayley, and Edward John Routh. Of Maxwell, Hopkins is reported to have said that he was the most extraordinary man he had met with in the whole course of his experience, that it seemed impossible for him to think wrongly on any physical subject, but that in analysis he was far more deficient. (Other contemporaries also testified to Maxwell's preference for geometrical over analytical methods.) This shrewd assessment was later borne out by several important formulas advanced by Maxwell that obtained correct results from faulty mathematical arguments.

In 1854 Maxwell was second wrangler and first Smith's prizeman (the Smith's prize is a prestigious competitive award for an essay that incorporates original research). He was elected to a fellowship at Trinity, but, because his father's health was deteriorating, he wished to return to Scotland. In 1856 he was appointed to the professorship of natural philosophy at Marischal College, Aberdeen, but before the appointment was announced his father died. This was a great personal loss, for Maxwell had had a close relationship with his father. In June 1858 Maxwell married Katherine Mary Dewar, daughter of the principal of Marischal College. The union was childless and was described by his biographer as a "married life . . . of unexampled devotion."

In 1860 the University of Aberdeen was formed by a merger between King's College and Marischal College, and Maxwell was declared redundant. He applied for a vacancy at the University of Edinburgh, but he was turned down in favour of his school friend Tait. He then was appointed to the professorship of natural philosophy at King's College, London.

The next five years were undoubtedly the most fruitful of his career. During this period his two classic papers on the electromagnetic field were published, and his demonstration of colour photography took place. He was elected to the Royal Society in 1861. His theoretical and experimental work on the viscosity of gases also was undertaken during these years and culminated in a lecture to the Royal Society in 1866. He supervised the experimental determination of electrical units for the British Association for the Advancement of Science, and this work in measurement and standardization led to the establishment of the National Physical Laboratory. He also measured the ratio of electromagnetic and electrostatic units of electricity and confirmed that it was in satisfactory agreement with the velocity of light as predicted by his theory.

Later life. In 1865 he resigned his professorship at King's College and retired to the family estate in Glenlair. He continued to visit London every spring and served as external examiner for the Mathematical Tripos (exams) at Cambridge. In the spring and early summer of 1867 he toured Italy. But most of his energy during this period was devoted to writing his famous treatise on electricity and magnetism.

It was Maxwell's research on electromagnetism that es-

Career

Education and early contributions

Research
on electro-
magnetism

established him among the great scientists of history. In the preface to his *Treatise on Electricity and Magnetism* (1873), the best exposition of his theory, Maxwell stated that his major task was to convert Faraday's physical ideas into mathematical form. In attempting to illustrate Faraday's law of induction (that a changing magnetic field gives rise to an induced electromagnetic field), Maxwell constructed a mechanical model. He found that the model gave rise to a corresponding "displacement current" in the dielectric medium, which could then be the seat of transverse waves. On calculating the velocity of these waves, he found that they were very close to the velocity of light. Maxwell concluded that he could "scarcely avoid the inference that light consists in the transverse undulations of the same medium which is the cause of electric and magnetic phenomena."

Maxwell's theory suggested that electromagnetic waves could be generated in a laboratory, a possibility first demonstrated by Heinrich Hertz in 1887, eight years after Maxwell's death. The resulting radio industry with its many applications thus has its origin in Maxwell's publications.

Maxwell's
other
contribu-
tions to
science

In addition to his electromagnetic theory, Maxwell made major contributions to other areas of physics. While still in his 20s, Maxwell demonstrated his mastery of classical physics by writing a prizewinning essay on Saturn's rings, in which he concluded that the rings must consist of masses of matter not mutually coherent—a conclusion that was corroborated more than 100 years later by the first Voyager space probe to reach Saturn.

The Maxwell relations of equality between different partial derivatives of thermodynamic functions are included in every standard textbook on thermodynamics (see THERMODYNAMICS, PRINCIPLES OF). Though Maxwell did not originate the modern kinetic theory of gases, he was the first to apply the methods of probability and statistics in describing the properties of an assembly of molecules. Thus he was able to demonstrate that the velocities of molecules in a gas, previously assumed to be equal, must follow a statistical distribution (known subsequently as the Maxwell-Boltzmann distribution law). In later papers Maxwell investigated the transport properties of gases—*i.e.*, the effect of changes in temperature and pressure on viscosity, thermal conductivity, and diffusion.

Maxwell was far from being an abstruse theoretician. He was skillful in the design of experimental apparatus, as was shown early in his career during his investigations of colour vision. He devised a colour top with adjustable sectors of tinted paper to test the three-colour hypothesis of Thomas Young and later invented a colour box that made it possible to conduct experiments with spectral colours rather than pigments. His investigations of the colour theory led him to conclude that a colour photography could be produced by photographing through filters of the three primary colours and then recombining the images. He demonstrated his supposition in a lecture to the Royal

Institution of Great Britain in 1861 by projecting through filters a colour photograph of a tartan ribbon that had been taken by this method.

In addition to these well-known contributions, a number of ideas that Maxwell put forward quite casually have since led to developments of great significance. The hypothetical intelligent being known as Maxwell's demon was a factor in the development of information theory. Maxwell's analytic treatment of speed governors is generally regarded as the founding paper on cybernetics, and his "equal areas" construction provided an essential constituent of the theory of fluids developed by Johannes Diederik van der Waals. His work in geometrical optics led to the discovery of the fish-eye lens. From the start of his career to its finish his papers are filled with novelty and interest. He also was a contributor to the ninth edition of *Encyclopædia Britannica*.

In 1871 Maxwell was elected to the new Cavendish professorship at Cambridge. He set about designing the Cavendish Laboratory and supervised its construction. Maxwell had few students, but they were of the highest calibre and included William D. Niven, Ambrose (later Sir Ambrose) Fleming, Richard Tetley Glazebrook, John Henry Poynting, and Arthur Schuster.

During the Easter term of 1879 Maxwell took ill on several occasions; he returned to Glenlair in June but his condition did not improve. He died after a short illness on Nov. 5, 1879. Maxwell received no public honours and was buried quietly in a small churchyard in the village of Parton, in Scotland.

BIBLIOGRAPHY. Maxwell's works include *Theory of Heat*, 3rd ed. (1872, reprinted 1970), and *A Treatise on Electricity and Magnetism*, 3rd ed., 2 vol. (1892, reissued 1954). Maxwell's original papers are collected in W.D. NIVEN (ed.), *The Scientific Papers of James Clerk Maxwell*, 2 vol. (1890, reissued 2 vol. in 1, 1965).

A standard biography is LEWIS CAMPBELL and WILLIAM GARNETT, *The Life of James Clerk Maxwell* (1882, reprinted 1969). Modern biographies include C.W.F. EVERITT, *James Clerk Maxwell: Physicist and Natural Philosopher* (1975), with a useful bibliography; IVAN TOLSTOY, *James Clerk Maxwell* (1981), for the general reader; and MARTIN GOLDMAN, *The Demon in the Aether* (1983). Commemorative publications include J.J. THOMSON *et al.*, *James Clerk Maxwell: A Commemorative Volume, 1831–1931* (1931), containing lectures given at Cambridge University; CYRIL DOMB (ed.), *Clerk Maxwell and Modern Science* (1963), lectures concerning his electromagnetic theory; and CYRIL DOMB, "James Clerk Maxwell: 100 Years Later." *Nature*, 282:235–239 (Nov. 15, 1979).

More details of his contributions to electromagnetism are in R.A.R. TRICKER, *The Contributions of Faraday and Maxwell to Electrical Science* (1966); EDMUND WHITTAKER, *A History of the Theories of Aether and Electricity*, rev. and enlarged ed., 2 vol. (1951–53, reprinted 1987); and JED Z. BUCHWALD, *From Maxwell to Microphysics: Aspects of Electromagnetic Theory in the Last Quarter of the Nineteenth Century* (1985), exploring Maxwellian theory and the transition into modern field theory. (Cy.Do./Ed.)

Cavendish
Laboratory

Measurement Systems

The concept of weights and measures, though it has come in modern times to include temperature, luminosity, pressure, and electric current, once consisted of only four basic measurements: mass (weight), volume (liquid or grain measure), distance, and area. The last three are, of course, closely related.

Basic to the whole idea of weights and measures are the concepts of uniformity, units, and standards. Uniformity, the essence of any system of weights and measures, requires accurate, reliable standards of mass and length and agreed-on units. A unit is the name of a quantity, such as kilogram or pound; a standard is the physical embodiment of a unit, such as the platinum-iridium cylinder kept by the International Bureau of Weights and Measures at Paris as the standard kilogram.

Two types of measurement systems are distinguished historically: an evolutionary system, such as the British Imperial, which grew more or less haphazardly out of custom, and a planned system, such as the present-day International System of Units (SI; *Système Internationale d'Unités*), in universal use by the world's scientific community and by most nations.

This article surveys the development of some of the world's major measurement systems and their distinctive characteristics. For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, section 723, and the *Index*.

This article is divided into the following sections:

Early units and standards	693
Ancient Mediterranean systems	
The ancient Chinese system	
Medieval systems	
The English and U.S. customary weights and measures systems	694
The English system	
The U.S. customary system	
The metric system of measurement	695
The development and establishment of the metric system	
The International System of Units	
Bibliography	697

EARLY UNITS AND STANDARDS

Ancient Mediterranean systems. Body measurements probably provided the most convenient bases for early linear measurements: early weight units may have derived casually from the use of certain containers or from calculations of what a person or animal could lift or haul.

The historical progression of units has followed a generally westward direction, the units of the ancient empires of the Middle East finding their way, mostly as a result of trade, to the Greek and then the Roman empires, thence to Gaul and Britain via Roman conquest.

The Egyptians. Although there is evidence that many early civilizations devised standards of measurement and some tools for measuring, the Egyptian cubit is generally recognized as having been the most ubiquitous standard of linear measurement in the very ancient world. Developed about 3000 BC, it was based on the length of the arm from the elbow to the extended fingertips and was standardized by a royal master cubit of black granite, against which all the cubit sticks in use in Egypt were measured at regular intervals.

The royal cubit (524 millimetres, or 20.62 inches) was subdivided in an extraordinarily complicated way. The basic subunit was the digit, doubtlessly a finger's breadth, of which there were 28 in the royal cubit. Four digits equaled a palm, five a hand. Twelve digits, or three palms, equaled a small span. Fourteen digits, or one-half a cubit, equaled

a large span. Sixteen digits, or four palms, made one *t'ser*. Twenty-four digits, or six palms, were a small cubit.

The digit was in turn subdivided. The 14th digit on a cubit stick was marked off into 16 equal parts. The next digit was divided into 15 parts, and so on, to the 28th digit, which was divided into 2 equal parts. Thus, measurement could be made to digit fractions with any denominator from 2 through 16. The smallest division, $\frac{1}{16}$ of a digit, was equal to $\frac{1}{448}$ part of a royal cubit.

The accuracy of the cubit stick is attested by the dimensions of the Great Pyramid of Giza; although thousands were employed in building it, its sides vary no more than 0.05 percent from the mean length of 230,364 metres (9,069.45 inches).

The Egyptians developed methods and instruments for measuring land at a very early date. The annual flood of the Nile River created a need for benchmarks and surveying techniques so that property boundaries could be readily reestablished when the water receded.

The Egyptian weight system appears to have been founded on a unit called the *kite*, with a decimal ratio, 10 *kites* equaling 1 *deben*, and 10 *debens* equaling 1 *sep*. Over the long span of Egyptian history, the weight of the *kite* varied from period to period, ranging all the way from 4.5 to 29.9 grams (0.16 to 1.05 ounce). About 3,400 different weights have been recovered from ancient Egypt, some in basic geometric shapes, others in human and animal forms.

Egyptian liquid measures, from large to small, were *ro*, *hin*, *hekat*, *khar*, and cubic cubit (0.14 cubic metre [37 U.S. gallons]).

The Babylonians. The earliest of all known weights is possibly the Babylonian mina, which in one surviving form weighed about 640 grams (23 ounces) and in another about 978 grams (34 ounces). Archaeologists have also found weights of 5 minas, in the shape of a duck, and a 30-mina weight in the form of a swan. The shekel, familiar from the Bible as a standard Hebrew coin and weight, was originally Babylonian. Most of the Babylonian weights and measures, carried in commerce throughout the Middle East, were gradually adopted by other countries. The basic Babylonian unit of length was the *kus* (about 530 millimetres, or 20.9 inches), also called the Babylonian cubit. The Babylonian *shusi*, defined as $\frac{1}{30}$ *kus*, was equal to 17.5 millimetres (0.69 inch). The Babylonian foot was $\frac{2}{3}$ *kus*.

The Babylonian liquid measure, *ka*, was the volume of a cube of one handbreadth (99 to 102 millimetres, or 3.9 to 4 inches). The cube, however, had to contain a weight of one great mina of water. The *ka* was a subdivision of two other units; 300 *ka* equaled 60 *gin* and 1 *gur*. The *gur* represented a volume of almost 303 litres (80 U.S. gallons).

The Hittites, Assyrians, Phoenicians, and Hebrews derived their systems generally from the Babylonians and Egyptians. Hebrew standards were based on the relationship between the mina, the talent (the basic unit), and the shekel. The sacred mina was equal to 60 shekels, and the sacred talent to 3,000 shekels, or 50 sacred minas. The Talmudic mina equaled 25 shekels; the Talmudic talent equaled 1,500 shekels, or 60 Talmudic minas.

The volumes of the several Hebrew standards of liquid measure are not definitely known; the *bat* may have contained about 37 litres (nearly 10 U.S. gallons); if so, the *log* equaled slightly more than 0.5 litre (0.13 U.S. gallon), and the *hin* slightly more than 6 litres (1.6 U.S. gallons). The Hebrew system was notable for the close relationship between dry and liquid volumetric measures; the liquid *kor* was the same size as the dry *homer*, and the liquid *bat* corresponded to the dry *'efa*.

Greeks and Romans. In the 1st millennium BC, commercial domination of the Mediterranean passed into the hands of the Greeks first and then the Romans. A basic

Early Greek units of measurement

Greek unit of length was the finger (19.3 millimetres, or 0.76 inch); 16 fingers equaled 30 centimetres (1 foot), and 24 fingers equaled 1 Olympic cubit. The coincidence with the Egyptian 24 digits equaling 1 small cubit suggests what is altogether probable on the basis of the commercial history of the era, that the Greeks derived their measures partly from the Egyptians and partly from the Babylonians, probably via the Phoenicians. The Greeks apparently used linear standards to establish their primary liquid measure, the *metrētēs*, equivalent to 39.4 litres (10.4 U.S. gallons). A basic Greek unit of weight was the talent (equal to 25.8 kilograms, or 56.9 pounds), obviously borrowed from Eastern neighbours.

The Romans, adapting the Greek system, subdivided the foot into 12 ounces, or inches (*unciae*), using the same word and the same subdivision, *unciae*, for $\frac{1}{12}$ pound (*libra*), which equaled 327.45 grams. The Romans made five feet equal to one pace, or double step, 1,000 of which made up the Roman mile (*mille passus*). The *sextarius* (0.53 litre, or 0.14 U.S. gallon) was the basic Roman unit of volume; it had several subdivisions and multiples, of which the largest, the amphora, borrowed from the Greeks, was 48 *sextarii* (25.5 litres, or 6.7 U.S. gallons).

The ancient Chinese system. Completely separated from the Mediterranean-European history of metrology is that of ancient China; yet the Chinese system exhibits all the principal characteristics of the Western. It employed parts of the body as a source of units—for example, the distance from the pulse to the base of the thumb. It was fundamentally chaotic in that there was no relationship between different types of units, such as those of length and those of volume. Finally, it was rich in variations. The *mou*, a unit of land measure, fluctuated from region to region from 0.08 to 0.13 hectare (0.2 to 0.3 acre). Variations were not limited to the geographic; a unit of length with the same name might be of one length for a carpenter, another for a mason, and still another for a tailor.

Shih huang-ti, who became the first emperor of China in 221 BC, is celebrated for, among other things, his unification of the regulations fixing the basic units. The basic weight, the *shih*, or *tan*, was fixed at about 60 kilograms (132 pounds); the two basic measurements, the *chih* and the *chang*, were set at about 25 centimetres (9.8 inches) and 3 metres (9.8 feet), respectively. A noteworthy characteristic of the Chinese system, and one that represented a substantial advantage over the Mediterranean systems, was its predilection for a decimal notation, as demonstrated by foot rulers dating back as far as the 6th century BC. Measuring instruments, too, were of a high order.

A unique characteristic of the Chinese system was its inclusion of an acoustic dimension. A standard vessel used for measuring grain and wine was defined not only as to weight but also as to pitch when struck; given a uniform shape and fixed weight, only a vessel of the proper volume would give the proper pitch. Thus the same word in old Chinese means "wine bowl," "grain measure," and "bell." Measures based on the length of a pitch pipe and its subdivision in terms of millet grains supplanted the old measurements based on the human body. The change brought a substantial increase in accuracy.

Medieval systems. Medieval Europe inherited the Roman system, with its Greek, Babylonian, and Egyptian roots. It soon proliferated through daily use and language variations into a great number of national and regional variants, with elements borrowed from Scandinavia and from the Arabs and original contributions growing out of the needs of medieval commerce.

A determined effort by Charlemagne to impose uniformity at the beginning of the 9th century was in vain; differing usages hardened. The great trade fairs, such as the Champagne Fairs of the 12th and 13th centuries, enforced rigid uniformity on merchants of all nationalities within the fairgrounds and had some effect on standardizing differences among regions, but the variations remained. A good example is the ell, the universal measure for wool cloth, the great trading staple of the Middle Ages. The ell of Champagne, two feet six inches, measured against an iron standard in the hands of the Keeper of the Fair, was accepted by Ypres and Ghent, both in modern Belgium;

by Arras, in modern France; and by the other great cloth-manufacturing cities of northwest Europe, even though their bolts varied in length. In several other parts of Europe, the ell itself varied, however.

The basic Roman unit of weight, the *libra*, acquired a Germanic name in parts of northern Europe but retained its Roman identity in the English abbreviation of pound as lb. Similarly the Roman mile survived, while the pace on which it was based vanished; it ceased to be a thousand of anything and instead became varying numbers of feet and yards, measures inherited from earlier northern Europe.

Medieval liquid measure was generally based on the *pinte*, or pint, which was approximately equal to the modern English quart; the quart was a medieval unit of dry measure, very close to its modern English equivalent in volume.

THE ENGLISH AND U.S. CUSTOMARY WEIGHTS AND MEASURES SYSTEMS

The English system. Out of the welter of medieval weights and measures emerged several national systems, reformed and reorganized from time to time; ultimately nearly all these were replaced by metric. In Britain and its American colonies, however, the ancient system survived.

By the time of Magna Carta (1215), abuses of weights and measures were so common that a clause was inserted in the charter to correct those on grain and wine. A few years later a royal ordinance entitled "Assize of Weights and Measures" defined a broad list of units and standards so successfully that it remained in force for nearly 600 years. A standard yard, "the Iron Yard of our Lord the King," was prescribed for the realm, divided into the traditional 3 feet, each of 12 inches, "neither more nor less." The perch (later the rod) was defined as 5.5 yards, or 16.5 feet. The inch was subdivided into 3 barley corns.

The furlong (a "furrow long") was eventually standardized as an eighth of a mile; the acre, probably from an old Anglo-Saxon word, as an area 4 rods wide by 40 long.

The influence of the Champagne Fairs may be seen in the separate English pounds for troy weight, probably from Troyes, one of the principal fair cities, and avoirdupois, the term used at the fairs for goods that had to be weighed—sugar, salt, alum, dyes, grain. The troy pound, for weighing gold and silver bullion and apothecaries' drugs, contained only 12 troy ounces.

A multiple of the English pound was the stone, which added a fresh element of confusion to the system by equaling neither 12 nor 16 but 14 pounds. The sacks of raw wool, which were medieval England's principal export, weighed 26 stones, or 364 pounds; huge standards, weighing 91 pounds, or one-fourth a sack, were employed in wool weighing. The sets of standards, which were sent out from London to the provincial towns, were usually of bronze or brass. Discrepancies somehow crept into the system, and in 1496, following a Parliamentary inquiry, new standards were made and sent out, a procedure repeated in 1588, under Elizabeth I.

No revision of law was found necessary for 200 years after Elizabeth's time, but several refinements and redefinitions were added. Edmund Gunter, a 17th-century mathematician, conceived the idea of taking the acre's breadth (4 perches, or 22 yards), calling it a chain, and dividing it into 100 links. In 1701 the corn bushel in dry measure was defined as "any round measure with a plain and even bottom, being 18.5 inches wide throughout and 8 inches deep." Similarly, in 1707 the wine gallon was defined as a round measure with an even bottom and containing 231 cubic inches; however, the ale gallon was retained at 282 cubic inches. There was also a corn gallon and an older, slightly smaller wine gallon.

A new Weights and Measures Act of 1824 sought to clear away some of the medieval tangle. A single gallon was decreed, defined as the volume occupied by "10 imperial pounds weight of distilled water weighed in air against brass weights with the water and the air at a temperature of 62 degrees of Fahrenheit's thermometer and with the barometer at 30 inches."

The same definition was reiterated in the Act of 1878, which redefined the yard: "the straight line or distance between the centres of two gold plugs or pins in the bronze

Troy and avoirdupois weights

Charlemagne's attempts to achieve standardization

bar . . . measured when the bar is at the temperature of sixty-two degrees of Fahrenheit's thermometer, and when it is supported by bronze rollers placed under it in such a manner as best to avoid flexure of the bar."

By an act of Parliament in 1963, all the English weights and measures were redefined in terms of the metric system, with a national changeover beginning two years later.

The U.S. customary system. In his first message to Congress in 1790, George Washington drew attention to the need for "uniformity in currency, weights and measures." Currency was settled in a decimal form, but the vast inertia of the English weights and measures system permeating industry and commerce and involving containers, measures, tools, and machines, as well as popular psychology, prevented the same approach, though it was advocated by Thomas Jefferson, from succeeding. In these very years the metric system was coming into being in France, and in 1821 Secretary of State John Quincy Adams, in a famous report to Congress, called the metric system "worthy of acceptance . . . beyond a question." Yet Adams admitted the impossibility of winning acceptance for it in the United States, until a future time "when the example of its benefits, long and practically enjoyed, shall acquire that ascendancy over the opinions of other nations which gives motion to the springs and direction to the wheels of the power."

Differences between the English and U.S. systems

Instead of adopting metric, the United States tried to bring its system into closer harmony with the English, from which various deviations had developed; for example, the United States still used "Queen Anne's gallon" of 231 cubic inches, which the British had discarded in 1824. Construction of standards was undertaken by the Office of Standard Weights and Measures, under the Treasury Department. The standard for the yard was one imported from London some years earlier, which guaranteed a close identity between the American and English yard; but Queen Anne's gallon was retained. The avoirdupois pound, at 7,000 grains, exactly corresponded with the British, as did the troy pound at 5,760 grains; however, the U.S. bushel, at 2,150.42 cubic inches, again deviated from the British. The U.S. bushel was derived from the "Winchester bushel," a surviving standard dating to the 15th century, which had been replaced in the British Act of 1824. It might be said that the U.S. gallon and bushel, smaller by about 17 percent and 3 percent, respectively, than the British, remain more truly medieval than their British counterparts.

At least the standards were fixed, however. From the mid-19th century, new states, as they were admitted to the Union, were presented with sets of standards. Late in the century, pressure grew to enlarge the role of the Office of Standard Weights and Measures, which, by Act of Congress effective July 1, 1901, became the National Bureau of Standards (since 1988 the National Institute of Standards and Technology), part of the Commerce Department. Its functions, as defined by the Act of 1901, included, besides the construction of physical standards and cooperation in establishment of standard practices, such activities as developing methods for testing materials and structures; carrying out research in engineering, physical science, and mathematics; and compilation and publication of general scientific and technical data. One of the first acts of the bureau was to sponsor a National Conference on Weights and Measures to coordinate standards among the states; one of the main functions of the annual conference became the updating of a model state law on weights and measures, which resulted in virtual uniformity in legislation.

Apart from this action, however, the U.S. government remained unique among major nations in refraining from exercising control at the national level. One noteworthy exception was the Metric Act of 1866, which permitted use of the metric system in the United States.

THE METRIC SYSTEM OF MEASUREMENT

The development and establishment of the metric system. One of the most significant results of the French Revolution was the establishment of the metric system of weights and measures.

European scientists had for many years discussed the desirability of a new, rational, and uniform system to replace the national and regional variants that made scientific communication difficult. The first proposal to closely approximate what eventually became the metric system was made as early as 1670. Gabriel Mouton, the vicar of St. Paul's Church in Lyon, suggested a length measure based on the arc of one minute of longitude, to be subdivided decimally. Mouton's proposal contained three of the major characteristics of the metric system: decimalization, rational prefixes, and the Earth's measurement as basis for a definition. Mouton's proposal was discussed, amended, criticized, and advocated for 120 years before the fall of the Bastille and the creation of the National Assembly made it a political possibility. In April of 1790 one of the foremost members of the assembly, Talleyrand, introduced the subject and launched a debate that resulted in a directive to the French Academy of Sciences to prepare a report. After several months' study, the academy recommended that the length of the meridian passing through Paris be determined from the North Pole to the Equator, that 1/10,000,000 of this distance be termed the metre and form the basis of a new decimal linear system, and, further, that a new unit of weight should be derived from the weight of a cubic metre of water. A list of prefixes for decimal multiples and submultiples was proposed. The National Assembly endorsed the report and directed that the necessary meridional measurements be taken.

The work of Gabriel Mouton

On June 19, 1791, a committee of 12 mathematicians, geodesists, and physicists met with Louis XVI, who gave his formal approval. The next day, the king attempted to escape from France, was arrested, returned to Paris, and was imprisoned; a year later, from his cell, he issued the proclamation that directed two engineers, Jean Delambre and Pierre Méchain, to perform the operations necessary to determine the length of the metre. The intervening time had been spent by the scientists and engineers in preliminary research; Delambre and Méchain now set to work to measure the distance on the meridian from Barcelona, Spain, to Dunkirk in northern France. The survey proved arduous; civil and foreign war so hampered the operation that it was not completed for six years. While Delambre and Méchain were struggling in the field, administrative details were being worked out in Paris. In 1793 a provisional metre was constructed from geodetic data already available. In 1795 the firm decision was taken to enact adoption of the metric system for France. The new law defined the length, mass, and capacity standards and listed the prefixes for multiples and submultiples. With the formal presentation to the assembly of the standard metre, as determined by Delambre and Méchain, the metric system became a fact in June 1799. The motto adopted for the new system was "For all people, for all time."

Establishment of the metric system

The standard metre was the Delambre-Méchain survey-derived "one ten-millionth part of a meridional quadrant of the earth." The gram, the basic unit of mass, was made equal to the mass of a cubic centimetre of pure water at the temperature of its maximum density (4° C [39.2° F]). A platinum cylinder known as the Kilogram of the Archives was declared the standard for 1,000 grams.

The litre was defined as the volume equivalent to the volume of a cube, each side of which had a length of 1 decimetre, or 10 centimetres.

The are was defined as the measure of area equal to a square 10 metres on a side. In practice the multiple hectare, 100 ares, became the principal unit of land measure.

The stère was defined as the unit of volume, equal to one cubic metre.

Names for multiples and submultiples of all units were made uniform, based on Greek prefixes.

The metric system's conquest of Europe was facilitated by the military successes of the French Revolution and Napoleon, but it required a long period of time to overcome the inertia of customary systems. Even in France Napoleon found it expedient to issue a decree permitting use of the old medieval system. Nonetheless, in the competition between the two systems existing side by side, the advantages of the metric proved decisive; in 1840 it was established as the legal monopoly in France, and from

Table 1: Elemental and Derived SI Units and Symbols

quantity	SI units		
	unit	formula	symbol
Elemental units			
Length	metre	—	m
Mass	kilogram	—	kg
Time	second	—	s
Electric current	ampere	—	A
Temperature	kelvin	—	K
Luminous intensity	candela	—	cd
Plane angle	radian	—	rad
Solid angle	steradian	—	sr
Derived units			
Acceleration	metre/second squared	m/s ²	
Area	square metre	m ²	
Capacitance	farad	A · s/V	F
Charge	coulomb	A · s	C
Density	kilogram/cubic metre	kg/m ³	
Electric field strength	volt/metre	V/m	
Energy	joule	N · m	J
Force	newton	kg · m/s ²	N
Frequency	hertz	s ⁻¹	Hz
Illumination	lux	lm/m ²	lx
Inductance	henry	V · s/A	H
Kinematic viscosity	square metre/second	m ² /s	
Luminance	candela/square metre	cd/m ²	
Luminous flux	lumen	cd · sr	lm
Magnetic field strength	ampere/metre	A/m	
Magnetic flux	weber	V · s	Wb
Magnetic flux density	tesla	Wb/m ²	T
Power	watt	J/s	W
Pressure	pascal (newton/square metre)	N/m ²	Pa
Resistance	ohm	V/A	Ω
Stress	pascal (newton/square metre)	N/m ²	Pa
Velocity	metre/second	m/s	
Viscosity	newton-second/square metre	N · s/m ²	
Voltage	volt	W/A	V
Volume	cubic metre	m ³	

that point forward its progress through the world has been steady, though it is worth observing that in many cases metric was adopted during the course of a political upheaval, just as in its original French beginning. Notable examples are Latin America, the Soviet Union, and China. In Japan the adoption of metric came about following the peaceful but far-reaching political changes associated with the Meiji Restoration of 1868.

In Britain, the Commonwealth nations, and the United States, the progress of metric has been discernible. The United States became a signatory to the Metric Convention

of 1875 and received copies of the International Prototype Metre and the International Prototype Kilogram in 1890. Three years later the Office of Weights and Measures announced that the prototype metre and kilogram would be regarded as fundamental standards from which the customary units, the yard and the pound, would be derived.

Throughout the 20th century, use of the metric system in various segments of commerce and industry increased spontaneously in Britain and the United States; it became almost universally employed in the scientific and medical professions. The automobile, electronics, chemical, and electric power industries have all adopted metric at least in part, as have such fields as optometry and photography. Legislative proposals to adopt metric generally have been made in the U.S. Congress and British Parliament. In 1968 Congress passed legislation calling for a program of investigation, research, and survey to determine the impact on the United States of increasing worldwide use of the metric system. The program concluded with a report to Congress in July 1971 that recommended: "on the basis of evidence marshalled in the U.S. metric study, this report [D.V. Simone, "A Metric America, A Decision whose Time has Come," *National Bureau of Standards Special Publication 345*] recommends that the United States change to the International Metric System through a coordinated national program over a period of ten years, at the end of which the nation will be predominantly metric." Parliament went further and established a long-range program of changeover.

In the meantime the metric system itself was replaced by the new International System of Units, which is fundamentally an expansion of 18th-century metric to incorporate scientific and technological developments of the 20th century.

The International System of Units. Just as the original conception of the metric system had grown out of the problems scientists encountered in dealing with the medieval system, so the new International System grew out of the problems a vastly enlarged scientific community faced in the proliferation of subsystems improvised to serve particular disciplines. At the same time, it had long been known that the original 18th-century standards were not accurate to the degree demanded by 20th-century scientific operations; new definitions were required. After lengthy discussion the 11th General Conference on Weights and Measures (11th CGPM), meeting in Paris, birthplace of

Increasing adoption of the metric system in the 20th century

Table 2: Common Equivalents and Conversion Factors for U.S. Customary and SI Systems

approximate common equivalents		conversions accurate within 10 parts per million	
1 inch	= 25 millimetres	inches × 25.4*	= millimetres
1 foot	= 0.3 metre	feet × 0.3048*	= metres
1 yard	= 0.9 metre	yards × 0.9144*	= metres
1 mile	= 1.6 kilometres	miles × 1.60934	= kilometres
1 square inch	= 6.5 square centimetres	square inches × 6.4516*	= square centimetres
1 square foot	= 0.09 square metre	square feet × 0.0929030	= square metres
1 square yard	= 0.8 square metre	square yards × 0.836127	= square metres
1 acre	= 0.4 hectare†	acres × 0.404686	= hectares
1 cubic inch	= 16 cubic centimetres	cubic inches × 16.3871	= cubic centimetres
1 cubic foot	= 0.03 cubic metre	cubic feet × 0.0283168	= cubic metres
1 cubic yard	= 0.8 cubic metre	cubic yards × 0.764555	= cubic metres
1 quart (liq)	= 1 litre†	quarts (liq) × 0.946353	= litres
1 gallon	= 0.004 cubic metre	gallons × 0.00378541	= cubic metres
1 ounce (avdp)	= 28 grams	ounces (avdp) × 28.3495	= grams
1 pound (avdp)	= 0.45 kilogram	pounds (avdp) × 0.453592	= kilograms
1 horsepower	= 0.75 kilowatt	horsepower × 0.745700	= kilowatts
1 millimetre	= 0.04 inch	millimetres × 0.0393701	= inches
1 metre	= 3.3 feet	metres × 3.28084	= feet
1 metre	= 1.1 yards	metres × 1.09361	= yards
1 kilometre	= 0.6 mile (statute)	kilometres × 0.621371	= miles (statute)
1 square centimetre	= 0.16 square inch	square centimetres × 0.155000	= square inches
1 square metre	= 11 square feet	square metres × 10.7639	= square feet
1 square metre	= 1.2 square yards	square metres × 1.19599	= square yards
1 hectare†	= 2.5 acres	hectares × 2.47105	= acres
1 cubic centimetre	= 0.06 cubic inch	cubic centimetres × 0.0610237	= cubic inches
1 cubic metre	= 35 cubic feet	cubic metres × 35.3147	= cubic feet
1 cubic metre	= 1.3 cubic yards	cubic metres × 1.30795	= cubic yards
1 litre†	= 1 quart (liq)	litres × 1.05669	= quarts (liq)
1 cubic metre	= 264 gallons	cubic metres × 264.172	= gallons
1 gram	= 0.035 ounce (avdp)	grams × 0.0352740	= ounces (avdp)
1 kilogram	= 2.2 pounds (avdp)	kilograms × 2.20462	= pounds (avdp)
1 kilowatt	= 1.3 horsepower	kilowatts × 1.34102	= horsepower

*Exact. †Common term not used in SI.
Source: National Bureau of Standards Wall Chart.

Table 3: Prefixes and Their Symbols in the SI System

multiples and submultiples	prefixes	symbols
10^{18}	exa	E
10^{15}	peta	P
10^{12}	tera	T
10^9	giga	G
10^6	mega	M
10^3	kilo	k
10^2	hecto	h
10	deca	da
10^{-1}	deci	d
10^{-2}	centi	c
10^{-3}	milli	m
10^{-6}	micro	μ
10^{-9}	nano	n
10^{-12}	pico	p
10^{-15}	femto	f
10^{-18}	atto	a

the metric system, in October 1960 formulated a new International System of Units (abbreviated SI). The SI was amended by subsequent convocations of the CGPM. The following base units have been adopted and defined:

Length: metre. Since 1983 the metre has been defined as the distance traveled by light in a vacuum in $1/299,792,458$ second.

Mass: kilogram. The standard for the unit of mass, the kilogram, is a cylinder of platinum-iridium alloy kept by the International Bureau of Weights and Measures, located in Sèvres, near Paris. A duplicate in the custody of the National Institute of Standards and Technology serves as the mass standard for the United States. (This is the only base unit still defined by an artifact.)

Time: second. The second is defined as the duration of 9,192,631,770 cycles of the radiation associated with a specified transition, or change in energy level, of the cesium-133 atom.

Electric current: ampere. The ampere is defined as the magnitude of the current that, when flowing through each of two long parallel wires separated by one metre in free space, results in a force between the two wires (due to their magnetic fields) of 2×10^{-7} newton (the newton is a unit of force equal to about 0.2 pound) for each metre of length.

Thermodynamic temperature: kelvin. The thermodynamic, or Kelvin, scale of temperature used in SI has its origin or zero point at absolute zero and has a fixed point at the triple point of water (the temperature and pressure at which ice, liquid water, and water vapour are in equilibrium), defined as 273.16 kelvins. The Celsius scale is derived from the Kelvin scale. The triple point is defined as 0.01° on the Celsius scale, which is approximately 32.02° on the Fahrenheit scale.

Amount of substance: mole. The mole is defined as the amount of substance containing the same number of chemical units (atoms, molecules, ions, electrons, or other

specified entities or groups of entities) as exactly 12 grams of carbon-12.

Light (luminous) intensity: candela. The candela is defined as the luminous intensity in a given direction of a source that emits monochromatic radiation at a frequency of 540×10^{12} hertz and that has a radiant intensity in the same direction of 1/683 watt per steradian (unit solid angle).

BIBLIOGRAPHY. The understanding and development of systems of physical measurements, from early elementary to sophisticated modern ones, are discussed in A.E. BERRIMAN, *Historical Metrology: A New Analysis of the Archaeological and the Historical Evidence Relating to Weights and Measures* (1953, reprinted 1969); B.D. ELLIS, *Basic Concepts of Measurement* (1966); BRUNO KISCH, *Scales and Weights: A Historical Outline* (1965); KEITH ELLIS, *Man and Measurement* (1973); H. ARTHUR KLEIN, *The World of Measurements: Masterpieces, Mysteries, and Muddles of Metrology* (1974; reissued 1988 as *The Science of Measurement: A Historical Survey*); OWEN BISHOP, *Yardsticks of the Universe* (1984); O.A.W. DILKE, *Mathematics and Measurement* (1987); and WITOLD KULA, *Measures and Men* (1986; originally published in Polish, 1970).

Advance toward modern systems of measurement is traced in ARTHUR E. KENNELLY, *Vestiges of Pre-Metric Weights and Measures Persisting in Metric-System Europe, 1926-1927* (1928); *Landmarks in Metrology—1983* (1983), a collection of papers written in the third quarter of the 20th century and published in connection with an international conference on the subject; REXMOND C. COCHRANE, *Measures for Progress: A History of the National Bureau of Standards* (1966, reprinted 1976), emphasizing the extent of NBS involvement in 20th-century scientific progress; C. DOUGLAS WOODWARD, *BSI—The Story of Standards* (1972), on the work of the British Standards Institution; LAL C. VERMAN, *Standardization, a New Discipline* (1973); and E.L. DELLOW, *Measuring and Testing in Science and Technology* (1970).

Development and maintenance of the U.S. customary system in particular are discussed in RALPH W. SMITH, *The Federal Basis for Weights and Measures: A Historical Review of Federal Legislative Effort, Statutes, and Administrative Action in the Field of Weights and Measures in the United States* (1958); L.J. CHISHOLM, *Units of Weight and Measure: International (Metric) and U.S. Customary* (1967, reprinted 1974), especially interesting for its tables of conversion equivalents; and LEWIS V. JUDSON, *Weights and Measures Standards of the United States: A Brief History* (1976).

Useful reference information, in dictionary or table form, is compiled in STEPHEN DRESNER, *Units of Measurement: An Encyclopaedic Dictionary of Units, Both Scientific and Popular, and the Quantities They Measure* (1971); JOHN L. FEIRER, *SI Metric Handbook* (1977); H.G. JERRARD and D.B. MCNEILL, *A Dictionary of Scientific Units: Including Dimensionless Numbers and Scales*, 5th ed. (1986); J.V. DRAZIL, *Quantities and Units of Measurement: A Dictionary and Handbook* (1983); and *The World Measurement Guide: Editorial Information Compiled by THE ECONOMIST*, 4th rev. ed. (1980).

Current research and development in the field is reported in such periodicals as *M & C: Measurements & Control* (bi-monthly, U.S.); *Measurement and Control* (monthly, Great Britain); *Measurement Science & Technology* (monthly); and *Measurement: Journal of the International Measurement Confederation* (quarterly). (L.J.C./Ed.)

Mecca and Medina

Mecca and Medina are the most sacred cities of Islām. Muḥammad, known to Muslims as the Prophet, was born in Mecca, and, according to Islāmīc tradition, it was while living there that he received the first divine revelations leading him to the foundation of Islām. It is toward the city's sacred shrine, the Ka'bah, that Muslims turn five times daily in prayer, and it is to this destination that all devout Muslims are obliged to attempt a pilgrimage, or hajj (Arabic: *hajj*), at least once in their lives. Medina is celebrated as the place in which Muḥammad founded the first Muslim community after his migration, or Hegira (Arabic: *hijrah*), from Mecca (AD

622) and from which he conquered all of Arabia in his war against the polytheists of Mecca. A lesser pilgrimage is made to Muḥammad's tomb, adjacent to the Prophet's mosque in Medina.

The two cities are situated in west-central Saudi Arabia in a region known as the Hejaz, and, separated by some 250 miles (400 kilometres), Medina lies due north of Mecca. The two have been traditionally identified with one another since the time of Muḥammad, and, because of their sacred features, only Muslims are allowed to enter their precincts.

This article is divided into the following sections:

Mecca 698

- Physical and human geography 698
 - The landscape
 - The city site
 - Climate
 - Plant and animal life
 - The city layout
 - Housing
 - The people
 - The economy
 - Industry
 - Transportation
 - Administration and social conditions
 - Government
 - Public utilities

- Education
- Health
- History 699
- Medina 700
 - Physical and human geography 700
 - The landscape
 - The city site
 - The city layout
 - The people
 - The economy
 - Agriculture
 - Industry
 - Transportation
 - History 701
 - Bibliography 701

MECCA

The holiest of Muslim cities, Mecca (Arabic: Makkah) is located in the Širāt Mountains inland from the Red Sea coast of Saudi Arabia. The city has an area of about 10 square miles (26 square kilometres).

A historic city closely tied to tradition, Mecca has benefited from the Saudi Arabian oil economy of the 20th century, although income from the hajj is still important. The city has undergone vast improvements. The area around the religious shrines has been cleared; the mosque has been enlarged; housing and sanitation have been improved; and transportation facilities have been enhanced. As a result, Mecca can accommodate the continually increasing number of pilgrims, or hajjis.

Physical and human geography

THE LANDSCAPE

The city site. Mecca is situated at an elevation of 909 feet (277 metres) above sea level in the dry beds of the Wādī Ibrāhīm and several of its short tributaries. It is surrounded by the Širāt Mountains, the peaks of which include Mount (Jabal) Ajyad, which rises to 1,332 feet, and Mount Abū Qubays, which attains 1,220 feet, to the east and Mount Qu'ayqān, which reaches 1,401 feet, to the west. Mount Hirā' rises to 2,080 feet on the northeast and contains a cave in which Muḥammad sought isolation and visions before he became a prophet. It was also in this cave that he received the first verse (*āyah*) of the holy Qur'ān. South of the city is Mount Thawr (2,490 feet), which contains the cave in which the Prophet secreted himself from his Meccan enemies during the Hegira to Medina, the event that marks the beginning of the Muslim calendar.

Entrance to the city is gained through four gaps in the surrounding mountains. The passes lead from the northeast to Minā, 'Arafāt, and Aṭ-Ṭā'if; from the northwest to Medina; from the west to Jiddah; and from the south to Yemen

(Šan'ā). The gaps have also defined the direction of the contemporary expansion of the city.

Climate. Because of its relatively low-lying location, Mecca is threatened by seasonal flash floods despite the low amount of annual precipitation. Less than five inches (130 millimetres) of rain falls during the year, mainly in the winter months. Temperatures are high throughout the year and in summer may reach 113° F (45° C).

Plant and animal life. Plant and animal life are scarce and consist of species that can withstand the high degree of aridity and heat. Natural vegetation is sparse and includes tamarisks and various types of acacia. Wild animals in the vicinity include wild cats, wolves, hyenas, foxes, mongooses, and kangaroo rats (jerboas).

The city layout. The city centres on the Ḥaram Mosque (the Great Mosque), in which are situated the Ka'bah and the sacred well of Zamzam. The compact built-up area around the mosque comprises the old city, which stretches to the north and southwest but is limited on the east and west by the nearby mountains. The main avenues are Al-Mudda'ah and Sūq al-Layl to the north of the mosque and As-Sūq as-Saghīr to the south.

Since World War II, Mecca has expanded along the roads through the mountain gaps to the north, northwest, and west. Among the modern residential areas are Al-'Azīziyah and Al-Faysaliyah along the road to Minā and Aṣ-Zāhir, Az-Zahra'ā, and Shāri' al-Manšūr along the roads to Jiddah and Medina. Expansion has been accompanied by the construction of new streets in the old city and by the transformation of Mecca into a modern city, with fountains, built since the 1950s, in its four main squares. The Ḥaram Mosque is magnificent in its size and architecture and has been embellished and enlarged on numerous occasions throughout the centuries, most recently in a massive expansion by the government of Saudi Arabia in the 1980s and '90s. The state-of-the-art complex, now multilevel, includes an advanced communication network, air condi-

Mecca's religious shrines

tioning, escalators, and a complex network of pedestrian routes and tunnels, in addition to numerous aesthetic and artistic accompaniments. The mosque can accommodate one million worshippers at a time. Houses near the mosque have been razed, and it is now surrounded by open spaces and wide streets, which can be crossed through underground walkways, built to ease traffic.

Housing. Mecca's houses are more compacted in the old city than in the modern residential areas. Traditional buildings of two or three stories are built of local rock. The villas in the modern areas are constructed of concrete. Slum conditions can still be found in various parts of the city; the slum inhabitants are mainly poor pilgrims who, unable to finance their return home, remained in Mecca after arriving either for the hajj or for a lesser pilgrimage known as the *'umrah*.

THE PEOPLE

The population density in Mecca is high. Most of the people are concentrated in the old city, while densities in the modern residential areas are the lowest in the city.

During the month of pilgrimage, the city is swollen with between one million and two million worshippers from other parts of the country and from other Muslim nations.

Entrance into Mecca is permitted only to followers of Islām. It is, however, one of the most cosmopolitan cities in the world, containing people from the various countries throughout the world. People of the same national origin tend to live together in certain parts of the city.

THE ECONOMY

Arable land and water are scarce, and food must be imported. Vegetables and fruits are brought in daily from the surrounding wadis, such as Wādī Fāṭimah, from the Aṭ-Ṭā'if area to the east-southeast, and from the southern agricultural areas, such as Bilād Ghāmid and Bilād Zahrān. Foodstuffs are imported from abroad mainly through the port of Jiddah, about 45 miles (70 kilometres) to the west on the Red Sea.

Industry. Industry is limited and includes the manufacture of textiles, furniture, and utensils. The overall urban economy is commercial and service-oriented.

Transportation. Transportation and facilities related to the hajj are the main services. Mecca has no airport or water or rail services. It is well served, however, by the Jiddah seaport and airport and by intercity truck, bus, and taxi services. A local bus system was established in 1979. Paved roads link Mecca with the main cities of Saudi Arabia and neighbouring countries.

Because of the improvement of services, the number of pilgrims has increased. This annual influx brings a good income to the city, but it also results in a temporary population of some two million or more, all of whom need accommodations, food, water, electricity, transportation, and medical services. To answer the problem of accommodations, the Saudi government has erected huge tent cities each year for the hajj, although sporadic fires in these camps have caused a number of deaths.

In accordance with the prescribed route, all pilgrims have to be transported from Mecca to 'Arafāt, a distance of nearly 12 miles, during the early morning of the ninth day of the month of Dhū al-Ḥijjah. During the night of the same day, they must travel to Minā, which is almost two miles from Mecca; after three days, all are returned to Mecca. This problem has been met by the construction of a good road network, an adequate supply of vehicles, and traffic control.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. The governor of the city is the emir of Makkah *minṭaqah idārīyah* (administrative district), who is responsible for the maintenance of law and order in both the city and the *minṭaqah*; he is appointed by the king and is immediately responsible to the minister of the interior. The municipal council is responsible for the functioning of the municipality; the council was formed after World War II and has 14 members, who are locally elected and are then approved by the minister of the interior. Mecca is also the capital of Makkah *minṭaqah idārīyah*, which includes the cities of Jiddah and Aṭ-Ṭā'if.

Public utilities. Mecca depends on the surrounding wadis for drinking water. The waters of Zubaydah Spring (Ayn Zubaydah), built in the 8th century, flow through tunnels from Wādī Nu'mān, about 20 miles to the southwest. Al-'Azīzīyah Spring sends its waters through pipelines from Wādī ash-Shāmīyah, about 60 miles to the northeast. Water and electricity have reached almost all houses. Electricity is generated at an oil-fueled power station located on the road to Medina.

Education. Free education is provided for both girls and boys from primary to university level. The Umm al-Qura University (founded 1979) is located in Mecca, as are two university colleges—the Madrasat Ahl al-Hadith and the Saudi Arabian Institute for Higher Education.

Health. Health services and medical care are free and adequate. Despite careful checking by officials, pilgrims may sometimes bring various illnesses, particularly cholera and cerebrospinal infections, into the city; the health services, however, have been able to keep such problems under control.

History

Ancient Mecca was an oasis on the old caravan trade route that linked the Mediterranean world with South Arabia, East Africa, and South Asia. The town was located about midway between Ma'rib in the south and Petra in the north, and it gradually developed by Roman and Byzantine times into an important trade and religious centre. It was known to Ptolemy as Macoraba.

According to Islāmic tradition, Abraham and Ishmael, his son by Hagar, built the Ka'bah as the house of God. The central point of pilgrimage in Mecca before the advent of Islām in the 7th century, the cube-shaped stone building has been destroyed and rebuilt several times. During pre-Islāmic times, the city was ruled by a series of Yemeni tribes. Under the Quraysh, it became a type of city-state, with strong commercial links to the rest of Arabia, Ethiopia, and Europe. Mecca became a place for trade, for pilgrimage, and for tribal gatherings.

The city gained its religious importance with the birth of Muḥammad about 570. The prophet was forced to flee from Mecca in 622, but he returned eight years later and took control of the city. He purged Mecca of idols, declared it a centre of Muslim pilgrimage, and dedicated it to God. Since then, the city has remained the major religious centre of Islām. As the ancient caravan route fell into decline, Mecca lost its commercial significance, and it has since lived mainly on the proceeds from the annual pilgrimages and the gifts of Muslim rulers. (A.S.A./Ed.)

Mecca was sacked by the Umayyad general al-Ḥajjāj ibn Yūsuf, and thereafter the city acknowledged the power of the Umayyad caliphate at Damascus and, following the eclipse of that dynasty, of the 'Abbāsīd caliphate of Baghdad. The city suffered great indignity at the hands of the Shī'ite Qarmatians in 930 when that sect's leader Ṭāhir Sulaymān pillaged Mecca and carried off the Black Stone from the Ka'bah. Beginning in the mid-10th century, the local city rulers were chosen from the sharifs, or descendants of Muḥammad, who retained a strong hold on the surrounding area while often paying homage to stronger political entities. The ability of the sharifs, originally moderate Shī'ites, to adapt to the changing political and religious climate ensured their preeminence in local affairs for the next 1,000 years. In 1269 Mecca came under the control of the Egyptian Mamlūk sultans. In 1517 dominion over the holy city passed to the Ottoman Empire, with its capital in Constantinople (now Istanbul). With the Ottoman collapse after World War I, control of Mecca was contested between the sharifs and the Āl Sa'ūd (the Sa'ūd family) of central Arabia, adherents to an austere, puritanical form of Islām known as Wahhābism. King Ibn Sa'ūd entered the city in 1925, and it became part of the Kingdom of Saudi Arabia and the capital of Makkah *minṭaqah idārīyah*.

Under Saudi rule, Wahhābism was enforced as the state credo, and the facilities for pilgrims were improved. Mecca underwent extensive economic development as Saudi Arabia's petroleum resources were exploited after World

Education services

The problems of the pilgrimage

The Ottoman period

War II, and the number of yearly pilgrims exploded. Despite lavish expenditures by the Saudi government to renovate the city and mosque area, in terms of both beauty and safety, the overwhelming crush of pilgrims each year has led to tragedy on several occasions, as in 1990, when nearly 1,500 pilgrims were trampled in a pedestrian tunnel, and in 1997, when several hundred died in a tent city fire and ensuing panic.

Political turmoil and violence have also often plagued the city. In 1979 a group of militants, mostly Saudi but many from other Islamic countries, seized the Haram Mosque and were evicted only with great loss of life after an assault by the Saudi National Guard. During the 1980s and '90s, Iranian pilgrims frequently engaged in political protests that led to clashes with Saudi police. Many deaths and injuries also ensued. (Ed.)

MEDINA

Medina (Arabic: Al-Madīnah) was known as Yathrib before the Prophet Muḥammad's residence there. The name was thereafter changed to Madīnat Rasūl Allāh ("City of the Messenger of God") or Al-Madīnah al-Munawwarah ("the Luminous City"). The city is situated in the Hejaz region of western Saudi Arabia, about 250 miles (400 kilometres) north of Mecca and 85 miles inland from the Red Sea. Medina is second only to Mecca as the holiest place of Islām.

Physical and human geography

THE LANDSCAPE

The city site. Medina lies some 2,000 feet (600 metres) above sea level on a fertile oasis. It is bounded on the east by an extensive lava field, part of which dates from a volcanic eruption in AD 1207. On the other three sides, the city is enclosed by arid hills belonging to the Hejaz mountain range. The highest of these hills is Mount Uḥūd, which rises to more than 2,000 feet above the oasis.

The city layout. Medina is a sacred area, and only Muslims are permitted to enter the city. The airport, however, lies just outside the sacred limits, and a good view of the city can be obtained by foreigners from aircraft landing at the airport.

In Turkish times, there was a small military landing ground at Sultanah, to the south near the garrison's barracks, but the area is now occupied by the king's palace and its extensive satellites. There too are the ruins of the tomb of 'Amr ibn al-Āṣ, the celebrated conqueror of Palestine and Egypt in early Islamic times. Other religious features of the oasis include the mosque of Qubā', the first in Islamic history; the Mosque of the Two Qiblahs, commemorating the change of the prayer direction from Jerusalem to Mecca, at Ar-Rimāh; the tomb of Ḥamzah, uncle of the Prophet, and of his companions who fell in the Battle of Uḥūd (625), in which the Prophet was wounded; and the cave in the flank of Uḥūd in which the Prophet took refuge on that occasion. Other mosques commemorate where he donned his armour for that battle; where he rested on the way thither, and where he unfurled his standard for the battle of the ditch (*khandaq*); and the ditch itself, dug around Medina by Muḥammad, in which the rubble of the great fire during the reign of Sultan Abdūlmecid I (1839–61) was dumped. All these spots are the object of pious visitation by all Muslims visiting Medina; they are forbidden to non-Muslims. In addition, the city is also the site of the Islamic University, established in 1961.

But the cynosure of all pilgrims is the Prophet's Mosque, which Muḥammad himself helped to build. Additions and improvements were undertaken by a succession of caliphs, and the chamber of the Prophet's wives was merged in the extension during the time of the Umayyad caliph al-Walid ibn 'Abd al-Malik. Fire twice damaged the mosque, first in 1256 and again in 1481, and its rebuilding was variously undertaken by devout rulers of several Islamic countries. Sultan Selim II (1566–74) decorated the interior of the mosque with mosaics overlaid with gold. Sultan Muḥammad built the dome in 1817 and in 1839 painted it green, this being the accepted colour of Islām. Sultan Abdūlmecid I initiated a project for the virtual reconstruction of the mosque in 1848 and completed it in 1860. A modern expansion was planned by King 'Abdul-'Aziz in 1948 and executed by King Sa'ūd in 1953–55. This renovation included the construction of a new northern court with its surrounding colonnades, all in the same style as the 19th-

century building but of concrete instead of stone from the neighbouring hills. The *Qafaṣ* (cage), to which female worshippers were formerly restricted, was dismantled, while, apart from minor repairs, the southern (main) part of the mosque remained intact. At its core, the mosque comprises the three ornamental iron structures representing the houses of the Prophet and containing, respectively (according to general consensus), the tomb of the Prophet himself under the great green dome, those of the first two caliphs (Abū Bakr and 'Umar), and that of the Prophet's daughter Fāṭimah. A specially adorned section of the pillared southern colonnade represents the palm grove (Ar-Rawḍah) in which the first simple mosque was built. In the 1980s and '90s, the government of Saudi Arabia lavished attention on the mosque, adding such accoutrements as state-of-the-art sliding domes, staircases, escalators, and an additional level for increased prayer space. The mosque is now equipped with an air-conditioning and shading system and an expansive logistic and transportation network. In addition, numerous aesthetic and structural embellishments were made; the mosque now has the capacity to hold roughly one million worshippers at a time.

Medina's modernization has not been so rapid as that of Jiddah, Riyadh, and other Saudi towns. In order to make room for new construction, the old city wall had to be completely dismantled, and that historic area had to merged with the now built-up pilgrim camping ground (Al-Manākh) and the Anbāriyah quarter, beyond the Abū-Jidā' torrent bed, which was formerly the commercial quarter and in which the Turks established the railway station and terminal yards. The foundations of the old city wall were found to be lower than the surface of accumulated silt and rubble; however, no attempt has been made to examine the excavations from an archaeological point of view. Nor has any archaeological work been done on the ruined sites of the old Jewish settlements, the largest of which was Yathrib (the Lathrippa or Iathrippa of Ptolemy and Stephanus Byzantius), which gave its name to the entire oasis until Islamic times. There are also several interesting mounds (*im*), besides the village of Al-Qurayzah, which would certainly produce historical data of interest. The Islamic cemetery of Al-Baqiyah was shorn of all the domes and ornamentation of the tombs of the saints at the time of the Wahhābī conquest of 1925; simple concrete graves in place of the old monuments and a circuit wall have been installed.

THE PEOPLE

The residents of Medina are Arabic-speaking Muslims, most of whom belong to the Sunnite branch of Islām. The city is one of the most populous in Saudi Arabia, and it is common for Muslims who make the pilgrimage to settle in the city. Farming and pottery making are important occupations.

THE ECONOMY

Agriculture. To supplement the income derived from accommodating pilgrims, Medina has an economy based on the cultivation of fruits, vegetables, and cereals. The city is especially well known for its date palms, the fruits of which are processed and packaged for export at a plant built in 1953.

Mechanical pumps for irrigation, in use since Turkish times early in the 20th century, have virtually replaced the old draw wells. Drinking water is supplied by an aqueduct from a spring at the southern end of the oasis. In addition

to the plentiful supply of subsoil water at no great depth, a number of important wadis meet in the vicinity of Medina and bring down torrents of water during the winter rains. Of these the most notable are the Wādī al-'Aqīq from the western mountains and a wadi coming down from the Aṭ-Ṭā'if area to the south.

Industry. Although Medina was known in early Islāmic times for metalworking, jewelry, and armoury, these industries were never large-scale, and most activity was connected with agricultural technology until the mid-20th century. Principal activities include automobile repair, brick and tile making, carpentry, and metalworking.

Transportation. From 1908 to 1916 Medina was connected with Damascus by the Hejaz railway, destroyed during World War I. Reconstruction of this railroad is studied periodically but has never taken place. Asphalt roads link the city with Jiddah, Mecca, and Yanbu' (Medina's port on the Red Sea); another road extends north through the Hejaz and connects the city to Jordan. Al-Jiladain airport nearby provides transportation to Saudi Arabian centres and has links to Jordan, Egypt, and Syria.

History

The earliest history of Medina is obscure, though it is known that there were Jewish settlers there in pre-Christian times. But the main influx of Jews is thought to have taken place as the result of their expulsion from Palestine by the Roman emperor Hadrian in about AD 135. It is probable that the Arab tribes of Aws and Khazraj were then in occupation of the oasis, but the Jews were the dominant factor in the population and development of the area by AD 400. On Sept. 20, 622, the arrival of the Prophet Muḥammad at Medina, in flight from Mecca, introduced a new chapter into the history of the oasis. Soon after the Hegira, tensions developed between the Jewish and Muslim communities, and the major Jewish groups were driven from the region. Medina became the administrative capital of the steadily expanding Islāmic state, a position it maintained until 661, when it was superseded in that role by Damascus, the capital of the Umayyad caliphs.

After the caliph's sack of the city in 683 for its fractious-

ness, the native emirs enjoyed a fluctuating measure of independence, interrupted by the aggressions of the sharifs of Mecca or controlled by the intermittent Egyptian protectorate.

The Ottomans, following their conquest of Egypt, held Medina after 1517 with a firmer hand, but their rule weakened and was almost nominal long before the Wahhābīs, an Islāmic revivalist group, first took the city in 1804. A Turko-Egyptian force retook it in 1812, and the Turks remained in effective control until the revival of the Wahhābī movement under Ibn Sa'ūd after 1912. Between 1904 and 1908 the Turks built the Hejaz railroad to Medina from Damascus in an attempt at strengthening the empire and ensuring Ottoman control over the hajj, the obligatory Muslim pilgrimage to the nearby holy city of Mecca. Turkish rule ceased during World War I, when Ḥusayn ibn 'Alī, the sharif of Mecca, revolted and put the railroad out of commission, with the assistance of the British officer T.E. Lawrence ("Lawrence of Arabia"). Ḥusayn later came into conflict with Ibn Sa'ūd, and in 1925 the city fell to the Sa'ūd dynasty. (J.B.GI./Ed.)

BIBLIOGRAPHY. Literature about Mecca is available mainly in Arabic. For pre-Islāmic and early Islāmic times, see M.A.A. AL-AZRAQI, *Akhbār Makkah*, written in the 9th century (1875, reprinted 1969); and A.I. AL-SHARIF, *Makkah wa-al-Madīnah* (1965). For the Middle Ages, see IBN JUBAYR, *The Travels of Ibn Jubayr* . . . , written in the 12th century, trans. by R.J.C. BROADHURST (1952); and IBN BATUTA, *Travels, A.D. 1325-1354*, written in the 14th century, trans. by H.A.R. GIBB (1958). Other English-language accounts include JOHN L. BURCKHARDT, *Travels in Arabia* (1829); JOHN F. KEANE, *Six Months in Meccah* (1881); C. SNOUCK HURGRONJE, *Mekka*, 2 vol. (1888-89); Eng. trans. of vol. 2, *Mekka in the Latter Part of the 19th Century: Daily Life, Customs and Learning* (1931); SIR RICHARD BURTON, *Personal Narrative of a Pilgrimage to El-Medīnah and Meccah*, 5th ed., 3 vol. (1906); ARTHUR WAVELL, *A Modern Pilgrim in Mecca and a Siege in Sanaa* (1912); ELDON RUTTER, *The Holy Cities of Arabia* (1928); J.B. PHILBY, *A Pilgrim in Arabia* (1946); DESMOND STEWART, *Mecca* (1980); JOHN SABINI, *Armies in the Sand: The Struggle for Mecca and Medina* (1981); and M.S. MAKKI, *Medina, Saudi Arabia: A Geographic Analysis of the City and Region* (1982). For statistics, see the *Statistical Yearbook* (annual), published by the Central Department of Statistics, Riyadh, Saudi Arabia.

The flight
from
Mecca

Mechanics: Energy, Forces, and Their Effects

Mechanics is the science concerned with the motion of bodies under the action of forces, including the special case in which a body remains at rest. Of first concern in the problem of motion are the forces that bodies exert on one another. This leads to the study of such topics as gravitation, electricity, and magnetism, according to the nature of the forces involved. Given the forces, one can seek the manner in which bodies move under the action of forces; this is the subject matter of mechanics proper.

Historically, mechanics was among the first of the exact sciences to be developed. Its internal beauty as a mathematical discipline and its early remarkable success in accounting in quantitative detail for the motions of the Moon, the Earth, and other planetary bodies had enormous influence on philosophical thought and provided impetus for the systematic development of science into the 20th century.

Mechanics may be divided into three branches: statics, which deals with forces acting on and in a body at rest; kinematics, which describes the possible motions of a body or system of bodies; and kinetics, which attempts to explain or predict the motion that will occur in a given situation. Alternatively, mechanics may be divided according to the kind of system studied. The simplest mechanical system is the particle, defined as a body so small that its shape and internal structure are of no consequence in the given problem. More complicated is the motion of a system of two or more particles that exert forces on one another and possibly undergo forces exerted by bodies outside of the system.

The principles of mechanics have been applied to three general realms of phenomena. The motions of such celestial bodies as stars, planets, and satellites can be predicted with great accuracy thousands of years before they occur. (The theory of relativity predicts some deviations from the motion according to classical, or Newtonian, mechanics; however, these are so small as to be observable only with very accurate techniques, except in problems involving all or a large portion of the detectable universe.) As the second realm, ordinary objects on Earth down to microscopic size (moving at speeds much lower than that of light) are properly described by classical mechanics without significant corrections. The engineer who designs bridges or aircraft may use the Newtonian laws of classical mechanics with confidence, even though the forces may be very complicated, and the beautifully simple and precise calculations of celestial mechanics usually cannot be duplicated. The third realm of phenomena comprises the behaviour of matter and electromagnetic radiation on the atomic and subatomic scale. Although there were some limited early successes in describing the behaviour of atoms in terms of classical mechanics, these phenomena are properly treated in quantum mechanics.

This article treats the fundamental parameters and concepts of classical mechanics and of certain key allied fields. It also considers those of quantum mechanics in some detail.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, sections 126–131, and the *Index*. (Ed.)

The article is divided into the following sections:

-
- | | |
|--|---|
| Classical mechanics 703 | History 735 |
| The origins and foundations of mechanics 703 | Basic principles 739 |
| History | Linear and angular momentum principles: |
| Fundamental concepts | stress and equations of motion |
| Motion of a particle in one dimension 706 | Geometry of deformation |
| Uniform motion | Stress-strain relations |
| Falling bodies and uniformly accelerated motion | Problems involving elastic response |
| Simple harmonic oscillations | Fluid mechanics 747 |
| Damped and forced oscillations | Basic properties of fluids 747 |
| Motion of a particle in two or more dimensions 709 | Hydrostatics 748 |
| Projectile motion | Hydrodynamics 750 |
| Motion of a pendulum | Bernoulli's law |
| Circular motion | Waves on shallow water |
| Circular orbits | Compressible flow in gases |
| Angular momentum and torque | Viscosity |
| Motion of a group of particles 713 | Navier-Stokes equation |
| Centre of mass | Potential flow |
| Conservation of momentum | Potential flow with circulation: vortex lines |
| Collisions | Waves on deep water |
| Relative motion | Boundary layers and separation |
| Coupled oscillators | Drag |
| Rigid bodies 717 | Lift |
| Statics | Turbulence |
| Rotation about a fixed axis | Convection |
| Rotation about a moving axis | Quantum mechanics 762 |
| Motion in a rotating frame | Historical basis of quantum theory 762 |
| Spinning tops and gyroscopes | Basic considerations |
| Analytic approaches 721 | Early developments |
| Configuration space | Planck's radiation law |
| The principle of virtual work | Einstein and the photoelectric effect |
| Lagrange's and Hamilton's equations | Bohr's theory of the atom |
| Celestial mechanics 723 | Scattering of X rays |
| Historical background | Broglie's wave hypothesis |
| Perturbations and problems of two bodies | Basic concepts and methods 763 |
| The three-body problem | Schrödinger's wave mechanics |
| The <i>n</i> -body problem | Electron spin and antiparticles |
| Tidal evolution | Identical particles and multielectron atoms |
| Relativistic mechanics 730 | Time-dependent Schrödinger equation |
| Development of the special theory of relativity | Tunneling |
| Relativistic space-time | Axiomatic approach |
| Relativistic momentum, mass, and energy | Incompatible observables |
| Mechanics of solids 734 | Heisenberg uncertainty principle |

Quantum electrodynamics	
The interpretation of quantum mechanics	768
The electron: wave or particle?	
Hidden variables	
Paradox of Einstein, Podolsky, and Rosen	
Measurement in quantum mechanics	

Applications of quantum mechanics	770
Decay of the K^0 meson	
Cesium clock	
A quantum voltage standard	
Bibliography	772

CLASSICAL MECHANICS

Classical mechanics deals with the motion of bodies under the influence of forces or with the equilibrium of bodies when all forces are balanced. The subject may be thought of as the elaboration and application of basic postulates first enunciated by Isaac Newton in his *Philosophiæ Naturalis Principia Mathematica* (1687), commonly known as the *Principia*. These postulates, called Newton's laws of motion, are set forth below. They may be used to predict with great precision a wide variety of phenomena ranging from the motion of individual particles to the interactions of highly complex systems. A variety of these applications are discussed in this article.

In the framework of modern physics, classical mechanics can be understood to be an approximation arising out of the more profound laws of quantum mechanics and the theory of relativity. However, that view of the subject's place greatly undervalues its importance in forming the context, language, and intuition of modern science and scientists. Our present-day view of the world and man's place in it is firmly rooted in classical mechanics. Moreover, many ideas and results of classical mechanics survive and play an important part in the new physics.

The central concepts in classical mechanics are force, mass, and motion. Neither force nor mass is very clearly defined by Newton, and both have been the subject of much philosophical speculation since Newton. Both of them are best known by their effects. Mass is a measure of the tendency of a body to resist changes in its state of motion. Forces, on the other hand, accelerate bodies, which is to say, they change the state of motion of bodies to which they are applied. The interplay of these effects is the principal theme of classical mechanics.

Although Newton's laws focus attention on force and mass, three other quantities take on special importance because their total amount never changes. These three quantities are energy, (linear) momentum, and angular momentum. Any one of these can be shifted from one body or system of bodies to another. In addition, energy may change form while associated with a single system, appearing as kinetic energy, the energy of motion; potential energy, the energy of position; heat, or internal energy, associated with the random motions of the atoms or molecules composing any real body; or any combination of the three. Nevertheless, the total energy, momentum, and angular momentum in the universe never changes. This fact is expressed in physics by saying that energy, momentum, and angular momentum are conserved. These three conservation laws arise out of Newton's laws, but Newton himself did not express them. They had to be discovered later.

It is a remarkable fact that, although Newton's laws are no longer considered to be fundamental, nor even exactly correct, the three conservation laws derived from Newton's laws—the conservation of energy, momentum, and angular momentum—remain exactly true even in quantum mechanics and relativity. In fact, in modern physics, force is no longer a central concept, and mass is only one of a number of attributes of matter. Energy, momentum, and angular momentum, however, still firmly hold centre stage. The continuing importance of these ideas inherited from classical mechanics may help to explain why this subject retains such great importance in science today.

The origins and foundations of mechanics

HISTORY

The discovery of classical mechanics was made necessary by the publication, in 1543, of the book *On the Revolutions*

of the Celestial Spheres by the Polish astronomer Nicolaus Copernicus. The book was about revolutions, real ones in the heavens, and it sparked the metaphorically named scientific revolution that culminated in Newton's *Principia* about 150 years later. The scientific revolution would change forever how people think about the universe.

In his book, Copernicus pointed out that the calculations needed to predict the positions of the planets in the night sky would be somewhat simplified if the Sun, rather than the Earth, were taken to be the centre of the universe (by which he meant what is now called the solar system). Among the many problems posed by Copernicus' book was an important and legitimate scientific question: if the Earth is hurtling through space and spinning on its axis as Copernicus' model prescribed, why is the motion not apparent?

To the casual observer, the Earth certainly seems to be solidly at rest. Scholarly thought about the universe in the centuries before Copernicus was largely dominated by the philosophy of Plato and Aristotle. According to Aristotelian science, the Earth was the centre of the universe. The four elements—earth, water, air, and fire—were naturally disposed in concentric spheres, with earth at the centre, surrounded respectively by water, air, and fire. Outside these were the crystal spheres on which the heavenly bodies rotated. Heavy, earthy objects fell because they sought their natural place. Smoke would rise through air, and bubbles through water for the same reason. These were natural motions. All other kinds of motion were violent motion and required a proximate cause. For example, an oxcart would not move without the help of an ox.

When Copernicus displaced the Earth from the centre of the universe, he tore the heart out of Aristotelian mechanics, but he did not suggest how it might be replaced. Thus, for those who wished to promote Copernicus' ideas, the question of why the motion of the Earth is not noticed took on a special urgency. Without suitable explanation, Copernicanism was a violation not only of Aristotelian philosophy but also of plain common sense.

The solution to the problem was discovered by the Italian mathematician and scientist Galileo Galilei. Inventing experimental physics as he went along, Galileo studied the motion of balls rolling on inclined planes. He noticed that, if a ball rolled down one plane and up another, it would seek to regain its initial height above the ground, regardless of the inclines of the two planes. That meant, he reasoned, that, if the second plane were not inclined at all but were horizontal instead, the ball, unable to regain its original height, would keep rolling forever. From this observation he deduced that bodies do not need a proximate cause to stay in motion. Instead, a body moving in the horizontal direction would tend to stay in motion unless something interfered with it. This is the reason that the Earth's motion is not apparent; the surface of the Earth and everything on and around it are always in motion together and therefore only seem to be at rest.

This observation, which was improved upon by the French philosopher and scientist René Descartes, who altered the concept to apply to motion in a straight line, would ultimately become Newton's first law, or the law of inertia. However, Galileo's experiments took him far beyond even this fundamental discovery. Timing the rate of descent of the balls (by means of precision water clocks and other ingenious contrivances) and imagining what would happen if experiments could be carried out in the absence of air resistance, he deduced that freely falling bodies would be uniformly accelerated at a rate independent of their mass. Moreover, he understood that the

Aristotelian
view of the
universe

Central
concepts
in classical
mechanics

motion of any projectile was the consequence of simultaneous and independent inertial motion in the horizontal direction and falling motion in the vertical direction. In his book *Dialogues Concerning the Two New Sciences* (1638), Galileo wrote.

It has been observed that missiles and projectiles describe a curved path of some sort; however, no one has pointed out the fact that this path is a parabola. But this and other facts, not few in number or less worth knowing, I have succeeded in proving . . .

Galileo's description of motion

Just as Galileo boasted, his studies would encompass many aspects of what is now known as classical mechanics, including not only discussions of the law of falling bodies and projectile motion but also an analysis of the pendulum, an example of harmonic motion. His studies fall into the branch of classical mechanics known as kinematics, or the description of motion. Although Galileo and others tried to formulate explanations of the causes of motion, the focus of the field termed dynamics, none would succeed before Newton.

Galileo's fame during his own lifetime rested not so much on his discoveries in mechanics as on his observations of the heavens, which he made with the newly invented telescope about 1610. What he saw there, particularly the moons of Jupiter, either prompted or confirmed his embrace of the Copernican system. At the time, Copernicus had few other followers in Europe. Among those few, however, was the brilliant German astronomer and mathematician Johannes Kepler.

Kepler devoted much of his scientific career to elucidating the Copernican system. Although Copernicus had put the Sun at the centre of the solar system, his astronomy was still rooted in the Platonic ideal of circular motion. Before Copernicus, astronomers had tried to account for the observed motions of heavenly bodies by imagining that they rotated on crystal spheres centred on the Earth. This picture worked well enough for the stars but not for the planets. To "save the appearances" (fit the observations) an elaborate system emerged of circular orbits, called epicycles, on top of circular orbits. This system of astronomy culminated with the *Almagest* of Ptolemy, who worked in Alexandria in the 2nd century AD. The Copernican innovation simplified the system somewhat, but Copernicus' astronomical tables were still based on circular orbits and epicycles. Kepler set out to find further simplifications that would help to establish the validity of the Copernican system.

In the course of his investigations, Kepler discovered the three laws of planetary motion that are still named for him. Kepler's first law says that the orbits of the planets are ellipses, with the Sun at one focus. This observation swept epicycles out of astronomy. His second law stated that, as the planet moved through its orbit, a line joining it to the Sun would sweep out equal areas in equal times. For Kepler, this law was merely a rule that helped him make precise calculations for his astronomical tables. Later, however, it would be understood to be a direct consequence of the law of conservation of angular momentum. Kepler's third law stated that the period of a planet's orbit depended only on its distance from the Sun. In particular, the square of the period is proportional to the cube of the semimajor axis of its elliptical orbit. This observation would suggest to Newton the inverse-square law of universal gravitational attraction.

By the middle of the 17th century, the work of Galileo, Kepler, Descartes, and others had set the stage for Newton's grand synthesis. Newton is thought to have made many of his great discoveries at the age of 23, when in 1665–66 he retreated from the University of Cambridge to his Lincolnshire home to escape from the bubonic plague. However, he chose not to publish his results until the *Principia* emerged 20 years later. In the *Principia*, Newton set out his basic postulates concerning force, mass, and motion. In addition to these, he introduced the universal force of gravity, which, acting instantaneously through space, attracted every bit of matter in the universe to every other bit of matter, with a strength proportional to their masses and inversely proportional to the square of the distance between them. These principles, taken together,

accounted not only for Kepler's three laws and Galileo's falling bodies and projectile motions but also for other phenomena, including the precession of the equinoxes, the oscillations of the pendulum, the speed of sound in air, and much more. The effect of Newton's *Principia* was to replace the by-then discredited Aristotelian worldview with a new, coherent view of the universe and how it worked. The way it worked is what is now referred to as classical mechanics.

FUNDAMENTAL CONCEPTS

Units and dimensions. Quantities have both dimensions, which are an expression of their fundamental nature, and units, which are chosen by convention to express magnitude or size. For example, a series of events have a certain duration in time. Time is the dimension of the duration. The duration might be expressed as 30 minutes or as half an hour. Minutes and hours are among the units in which time may be expressed. One can compare quantities of the same dimensions, even if they are expressed in different units (an hour is longer than a minute). Quantities of different dimensions cannot be compared with one another.

The fundamental dimensions used in mechanics are time, mass, and length. Symbolically, these are written as t , m , and l , respectively. The study of electromagnetism adds an additional fundamental dimension, electric charge, or q . Other quantities have dimensions compounded of these. For example, speed has the dimensions distance divided by time, which can be written as l/t , and volume has the dimensions distance cubed, or l^3 . Some quantities, such as temperature, have units but are not compounded of fundamental dimensions.

There are also important dimensionless numbers in nature, such as the number $\pi = 3.14159 \dots$. Dimensionless numbers may be constructed as ratios of quantities having the same dimension. Thus, the number π is the ratio of the circumference of a circle (a length) to its diameter (another length). Dimensionless numbers have the advantage that they are always the same, regardless of what set of units is being used.

Governments have traditionally been responsible for establishing and enforcing standard units for the sake of orderly commerce, navigation, science, and, of course, taxation. Today all such units are established by international treaty, revised every few years in light of scientific findings. The units used for most scientific measurements are those designated the International System of Units (Système International d'Unités), or SI for short. They are based on the metric system, first adopted officially by France in 1795. Other units, such as those of the British engineering system, are still in use in some places, but these are now defined in terms of the SI units.

The fundamental unit of length is the metre. A metre used to be defined as the distance between two scratch marks on a metal bar kept in Paris, but it is now much more precisely defined as the distance that light travels in a certain time interval ($1/299,792,458$ of a second). By contrast, in the British system, units of length have a clear human bias: the foot, the inch (the first joint of the thumb), the yard (distance from nose to outstretched fingertip), and the mile (one thousand standard paces of a Roman legion). Each of these is today defined as some fraction or multiple of a metre (one yard is nearly equal to one metre). In the SI or the metric system, lengths are expressed as decimal fractions or multiples of a metre (a millimetre = one-thousandth of a metre; a centimetre = one-hundredth of a metre; a kilometre = one thousand metres).

Times longer than one second are expressed in the units seconds, minutes, hours, days, weeks, and years. Times shorter than one second are expressed as decimal fractions (a millisecond = one-thousandth of a second, a microsecond = one-millionth of a second, and so on). The fundamental unit of time (*i.e.*, the definition of one second) is today based on the intrinsic properties of certain kinds of atoms (an excitation frequency of the isotope cesium-133).

Units of mass are also defined in a way that is technically sound, but in common usage they are the subject of some

Units of length

Great contributions of Newton

confusion because they are easily confused with units of weight, which is a different physical quantity. The weight of an object is the consequence of the Earth's gravity operating on its mass. Thus, the mass of a given object is the same everywhere, but its weight varies slightly if it is moved about the surface of the Earth, and it would change a great deal if it were moved to the surface of another planet. Also, weight and mass do not have the same dimensions (weight has the dimensions $m/l/t^2$). The Constitution of the United States, which calls on the government to establish uniform "weights and measures," is oblivious to this distinction, as are merchants the world over, who measure the weight of bread or produce but sell it in units of kilograms, the SI unit of mass. (The kilogram is equal to 1,000 grams; 1 gram is the mass of 1 cubic centimetre of water—under appropriate conditions of temperature and pressure.)

Vectors. The equations of mechanics are typically written in terms of Cartesian coordinates. At a certain time t , the position of a particle may be specified by giving its coordinates $x(t)$, $y(t)$, and $z(t)$ in a particular Cartesian frame of reference. However, a different observer of the same particle might choose a differently oriented set of mutually perpendicular axes, say, x' , y' , and z' . The motion of the particle is then described by the first observer in terms of the rate of change of $x(t)$, $y(t)$, and $z(t)$, while the second observer would discuss the rates of change of $x'(t)$, $y'(t)$, and $z'(t)$. That is, both observers see the same particle executing the same motion and obeying the same laws, but they describe the situation with different equations. This awkward situation may be avoided by means of a mathematical construction called a vector. Although vectors are mathematically simple and extremely useful in discussing mechanics, they were not developed in their modern form until late in the 19th century, when J. Willard Gibbs and Oliver Heaviside (of the United States and Britain, respectively) each applied vector analysis in order to help express the new laws of electromagnetism proposed by James Clerk Maxwell.

A vector is a quantity that has both magnitude and direction. It is typically represented symbolically by an arrow in the proper direction, whose length is proportional to the magnitude of the vector. Although a vector has magnitude and direction, it does not have position. A vector is not altered if it is displaced parallel to itself as long as its length is not changed.

By contrast to a vector, an ordinary quantity having magnitude but not direction is known as a scalar. In printed works vectors are often represented by boldface letters such as \mathbf{A} or \mathbf{X} , and scalars are represented by lightface letters, A or X . The magnitude of a vector, denoted $|\mathbf{A}|$, is itself a scalar—i.e., $|\mathbf{A}| = A$.

Because vectors are different from ordinary (i.e., scalar) quantities, all mathematical operations involving vectors must be carefully defined. Addition, subtraction, three kinds of multiplication, and differentiation will be discussed here. There is no mathematical operation that corresponds to division by a vector.

If vector \mathbf{A} is added to vector \mathbf{B} , the result is another vector, \mathbf{C} , written $\mathbf{A} + \mathbf{B} = \mathbf{C}$. The operation is performed by displacing \mathbf{B} so that it begins where \mathbf{A} ends, as shown in Figure 1A. \mathbf{C} is then the vector that starts where \mathbf{A} begins and ends where \mathbf{B} ends.

Vector addition is defined to have the (nontrivial) property $\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$. There do exist quantities having magnitude and direction that do not obey this requirement. An example is finite rotations in space. Two finite rotations of a body about different axes do not necessarily result in the same orientation if performed in the opposite order.

Vector subtraction is defined by $\mathbf{A} - \mathbf{B} = \mathbf{A} + (-\mathbf{B})$, where the vector $-\mathbf{B}$ has the same magnitude as \mathbf{B} but the opposite direction. The idea is illustrated in Figure 1B.

A vector may be multiplied by a scalar. Thus, for example, the vector $2\mathbf{A}$ has the same direction as \mathbf{A} but is twice as long. If the scalar has dimensions, the resulting vector still has the same direction as the original one, but the two cannot be compared in magnitude. For example, a particle moving with constant velocity \mathbf{v} suffers a displacement \mathbf{s} in time t given by $\mathbf{s} = \mathbf{v}t$. The vector \mathbf{v} has been multiplied

by the scalar t to give a new vector, \mathbf{s} , which has the same direction as \mathbf{v} but cannot be compared to \mathbf{v} in magnitude (a displacement of one metre is neither bigger nor smaller than a velocity of one metre per second). This is a typical example of a phenomenon that might be represented by different equations in differently oriented Cartesian coordinate systems but that has a single vector equation (for all observers not moving with respect to one another).

The dot product (also known as the scalar product, or sometimes the inner product) is an operation that combines two vectors to form a scalar. The operation is written $\mathbf{A} \cdot \mathbf{B}$. If θ is the (smaller) angle between \mathbf{A} and \mathbf{B} , then the result of the operation is $\mathbf{A} \cdot \mathbf{B} = AB \cos \theta$. The dot product measures the extent to which two vectors are parallel. It may be thought of as multiplying the magnitude of one vector (either one) by the projection of the other upon it, as shown in Figure 1C. If the two vectors are perpendicular, the dot product is zero.

The cross product (also known as the vector product) combines two vectors to form another vector, perpendicular to the plane of the original vectors. The operation is written $\mathbf{A} \times \mathbf{B}$. If θ is the (smaller) angle between \mathbf{A} and \mathbf{B} , then $|\mathbf{A} \times \mathbf{B}| = AB \sin \theta$. The direction of $\mathbf{A} \times \mathbf{B}$ is given

From R.P. Olenick, T.M. Apostol, and D.L. Goodstein, *The Mechanical Universe Introduction to Mechanics and Heat* (1985), Cambridge University Press

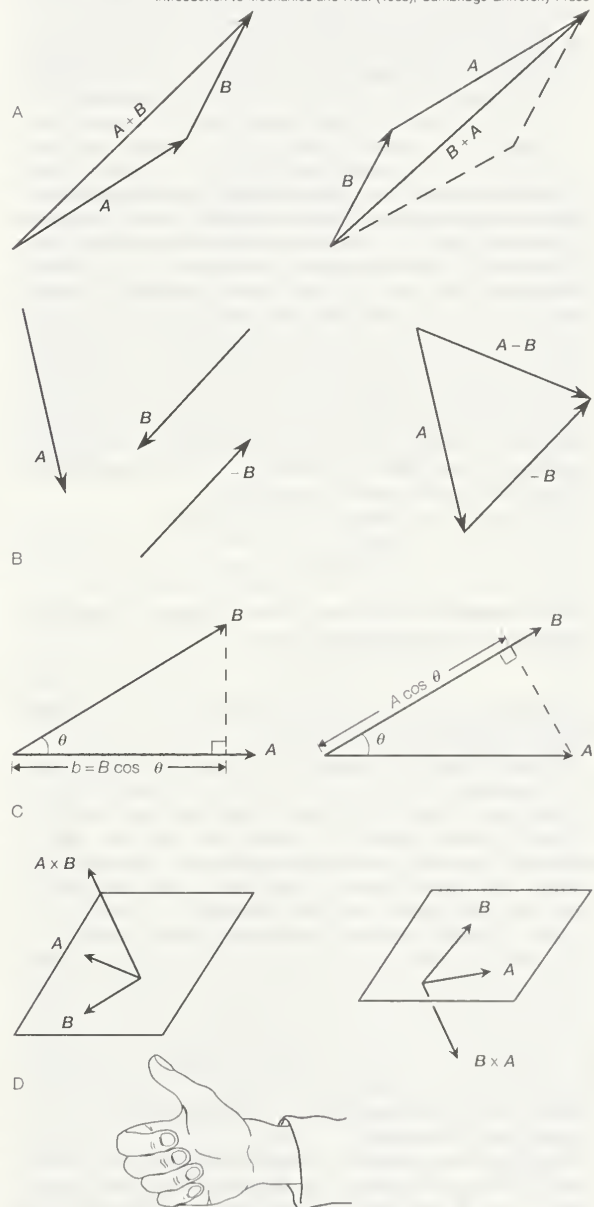


Figure 1: (A, left) The vector sum $\mathbf{A} + \mathbf{B}$ and (right) the vector sum $\mathbf{B} + \mathbf{A}$. (B, left) The vectors \mathbf{A} , \mathbf{B} , and $-\mathbf{B}$ and (right) the vector difference $\mathbf{A} - \mathbf{B}$. (C, left) $B \cos \theta$ is the component of \mathbf{B} along \mathbf{A} , and (right) $A \cos \theta$ is the component of \mathbf{A} along \mathbf{B} . (D, left) The right-hand rule used to find the direction of $\mathbf{A} \times \mathbf{B}$ and (right) the right-hand rule used to find the direction of $\mathbf{B} \times \mathbf{A}$.

Mathematical operations of vectors

Right-hand rule by the right-hand rule: if the fingers of the right hand are made to rotate from A through θ to B , the thumb points in the direction of $A \times B$, as shown in Figure 1D. The cross product is zero if the two vectors are parallel, and it is maximum in magnitude if they are perpendicular.

The derivative, or rate of change, of a vector is defined in perfect analogy to the derivative of a scalar: if the vector A changes with time t , then

$$\frac{dA}{dt} = \lim_{\Delta t \rightarrow 0} \frac{A(t + \Delta t) - A(t)}{\Delta t} \quad (1)$$

Before going to the limit on the right-hand side of equation (1), the operations described are vector subtraction [$A(t + \Delta t) - A(t)$] and scalar multiplication (by $1/\Delta t$). The result, dA/dt , is therefore itself a vector. Notice that, as shown in Figure 1B, the difference between two vectors, in this case $A(t + \Delta t) - A(t)$, may be in quite a different direction than either of the vectors from which it is formed, here $A(t + \Delta t)$ and $A(t)$. As a result, dA/dt may be in a different direction than $A(t)$.

Newton's laws of motion and equilibrium. In his *Principia*, Newton reduced the basic principles of mechanics to three laws:

1. Every body continues in its state of rest or of uniform motion in a straight line, unless it is compelled to change that state by forces impressed upon it.
2. The change of motion of an object is proportional to the force impressed and is made in the direction of the straight line in which the force is impressed.
3. To every action there is always opposed an equal reaction; or, the mutual actions of two bodies upon each other are always equal and directed to contrary parts.

Newton's first law is a restatement of the principle of inertia, proposed earlier by Galileo and perfected by Descartes. The second law is the most important of the three: it may be understood very nearly to summarize all of classical mechanics. Newton used the word "motion" to mean what is today called momentum—that is, the product of mass and velocity, or $p = mv$, where p is the momentum, m the mass, and v the velocity of a body. The second law may then be written in the form of the equation $F = dp/dt$, where F is the force, the time derivative expresses Newton's "change of motion," and the vector form of the equation assures that the change is in the same direction as the force, as the second law requires.

For a body whose mass does not change,

$$\frac{dp}{dt} = m \frac{dv}{dt} = ma,$$

where a is the acceleration. Thus, Newton's second law may be put in the following form:

$$F = ma. \quad (2)$$

It is probably fair to say that equation (2) is the most famous equation in all of physics.

Newton's third law assures that when two bodies interact, regardless of the nature of the interaction, they do not produce a net force acting on the two-body system as a whole. Instead, there is an action and reaction pair of equal and opposite forces, each acting on a different body (action and reaction forces never act on the same body). The third law applies whether the bodies in question are at rest, in uniform motion, or in accelerated motion.

If a body has a net force acting on it, it undergoes accelerated motion in accordance with the second law. If there is no net force acting on a body, either because there are no forces at all or because all forces are precisely balanced by contrary forces, the body does not accelerate and may be said to be in equilibrium. Conversely, a body that is observed not to be accelerated may be deduced to have no net force acting on it.

Consider, for example, a massive object resting on a table. The object is known to be acted on by the gravitational force of the Earth; if the table were removed, the object would fall. It follows therefore from the fact that the object does not fall that the table exerts an upward force on the object, equal and opposite to the downward force of gravity. This upward force is not a mere physicist's bookkeeping device but rather a real physical force.

The table's surface is slightly deformed by the weight of the object, causing the surface to exert a force analogous to that exerted by a coiled spring.

It is useful to recall the following distinction: the massive object exerts a downward force on the table that is equal and opposite to the upward force exerted by the table (owing to its deformation) on the object. These two forces are an action and reaction pair operating on different bodies (one on the table, the other on the object) as required by Newton's third law. On the other hand, the upward force exerted on the object by the table is balanced by a downward force exerted on the object by the Earth's gravity. These two equal and opposite forces, acting on the same body, are not related to or by Newton's third law, but they do produce the equilibrium immobile state of the body.

Motion of a particle in one dimension

UNIFORM MOTION

According to Newton's first law (also known as the principle of inertia), a body with no net force acting on it will either remain at rest or continue to move with uniform speed in a straight line, according to its initial condition of motion. In fact, in classical Newtonian mechanics, there is no important distinction between rest and uniform motion in a straight line; they may be regarded as the same state of motion seen by different observers, one moving at the same velocity as the particle, the other moving at constant velocity with respect to the particle.

Although the principle of inertia is the starting point and the fundamental assumption of classical mechanics, it is less than intuitively obvious to the untrained eye. In Aristotelian mechanics, and in ordinary experience, objects that are not being pushed tend to come to rest. The law of inertia was deduced by Galileo from his experiments with balls rolling down inclined planes, described above.

For Galileo, the principle of inertia was fundamental to his central scientific task: he had to explain how it is possible that if the Earth is really spinning on its axis and orbiting the Sun we do not sense that motion. The principle of inertia helps to provide the answer: Since we are in motion together with the Earth, and our natural tendency is to retain that motion, the Earth appears to us to be at rest. Thus, the principle of inertia, far from being a statement of the obvious, was once a central issue of scientific contention. By the time Newton had sorted out all the details, it was possible to account accurately for the small deviations from this picture caused by the fact that the motion of the Earth's surface is not uniform motion in a straight line (the effects of rotational motion are discussed below). In the Newtonian formulation, the common observation that bodies that are not pushed tend to come to rest is attributed to the fact that they have unbalanced forces acting on them, such as friction and air resistance.

As has already been stated, a body in motion may be said to have momentum equal to the product of its mass and its velocity. It also has a kind of energy that is due entirely to its motion, called kinetic energy. The kinetic energy of a body of mass m in motion with velocity v is given by

$$K = \frac{1}{2}mv^2. \quad (3)$$

FALLING BODIES AND UNIFORMLY ACCELERATED MOTION

During the 14th century, the French scholar Nicole Oresme studied the mathematical properties of uniformly accelerated motion. He had little interest in whether that kind of motion could be observed in the realm of actual human existence, but he did discover that, if a particle is uniformly accelerated, its speed increases in direct proportion to time, and the distance it traverses is proportional to the square of the time spent accelerating. Two centuries later, Galileo repeated these same mathematical discoveries (perhaps independently) and, just as important, determined that this kind of motion is actually executed by balls rolling down inclined planes. As the incline of the plane increases, the acceleration increases, but the

Significance of the second law of motion

Principle of inertia

motion continues to be uniformly accelerated. From this observation, Galileo deduced that a body falling freely in the vertical direction would also have uniform acceleration. Even more remarkably, he demonstrated that, in the absence of air resistance, all bodies would fall with the same constant acceleration regardless of their mass. If the constant acceleration of any body dropped near the surface of the Earth is expressed as g , the behaviour of a body dropped from rest at height z_0 and time $t = 0$ may be summarized by the following equations:

$$z = z_0 - \frac{1}{2}gt^2, \quad (4)$$

$$v = gt, \quad (5)$$

$$a = g, \quad (6)$$

where z is the height of the body above the surface, v is its speed, and a is its acceleration. These equations of motion hold true until the body actually strikes the surface. The value of g is approximately 9.8 metres per second squared (m/s^2).

A body of mass m at a height z_0 above the surface may be said to possess a kind of energy purely by virtue of its position. This kind of energy (energy of position) is called potential energy. The gravitational potential energy is given by

$$U = mgz_0. \quad (7)$$

Technically, it is more correct to say that this potential energy is a property of the Earth-body system rather than a property of the body itself, but this pedantic distinction can be ignored.

As the body falls to height z less than z_0 , its potential energy U converts to kinetic energy $K = \frac{1}{2}mv^2$. Thus, the speed v of the body at any height z is given by solving the equation

$$\frac{1}{2}mv^2 + mgz = mgz_0. \quad (8)$$

Equation (8) is an expression of the law of conservation of energy. It says that the sum of kinetic energy, $\frac{1}{2}mv^2$, and potential energy, mgz , at any point during the fall, is equal to the total initial energy, mgz_0 , before the fall began. Exactly the same dependence of speed on height could be deduced from the kinematic equations (4), (5), and (6) above.

In order to reach the initial height z_0 , the body had to be given its initial potential energy by some external agency, such as a person lifting it. The process by which a body or a system obtains mechanical energy from outside of itself is called work. The increase of the energy of the body is equal to the work done on it. Work is equal to force times distance.

The force exerted by the Earth's gravity on a body of mass m may be deduced from the observation that the body, if released, will fall with acceleration g . Since force is equal to mass times acceleration, the force of gravity is given by $F = mg$. To lift the body to height z_0 , an equal and opposite (*i.e.*, upward) force must be exerted through a distance z_0 . Thus, the work done is

$$W = Fz_0 = mgz_0, \quad (9)$$

which is equal to the potential energy that results.

If work is done by applying a force to a body that is not being acted upon by an opposing force, the body is accelerated. In this case, the work endows the body with kinetic energy rather than potential energy. The energy that the body gains is equal to the work done on it in either case. It should be noted that work, potential energy, and kinetic energy, all being aspects of the same quantity, must all have the dimensions ml^2/t^2 .

SIMPLE HARMONIC OSCILLATIONS

Consider a mass m held in an equilibrium position by springs, as shown in Figure 2A. The mass may be perturbed by displacing it to the right or left. If x is the displacement of the mass from equilibrium (Figure 2B), the springs exert a force F proportional to x , such that

$$F = -kx, \quad (10)$$

where k is a constant that depends on the stiffness of the springs. Equation (10) is called Hooke's law, and the force is called the spring force. If x is positive (displacement to the right), the resulting force is negative (to the left), and vice versa. In other words, the spring force always acts so as to restore mass back toward its equilibrium position.

Hooke's law

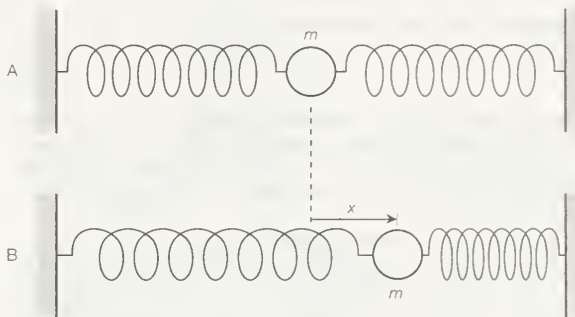


Figure 2: (A) A mass m held in equilibrium by springs. (B) A mass m displaced a distance x .

Moreover, the force will produce an acceleration along the x direction given by $a = d^2x/dt^2$. Thus, Newton's second law, $F = ma$, is applied to this case by substituting $-kx$ for F and d^2x/dt^2 for a , giving $-kx = m(d^2x/dt^2)$. Transposing and dividing by m yields the equation

$$a = \frac{d^2x}{dt^2} = -\frac{k}{m}x. \quad (11)$$

Equation (11) gives the derivative—in this case the second derivative—of a quantity x in terms of the quantity itself. Such an equation is called a differential equation, meaning an equation containing derivatives. Much of the ordinary, day-to-day work of theoretical physics consists of solving differential equations. The question is, given equation (11), how does x depend on time?

Solution of differential equations

The answer is suggested by experience. If the mass is displaced and released, it will oscillate back and forth about its equilibrium position. That is, x should be an oscillating function of t , such as a sine wave or a cosine wave. For example, x might obey a behaviour such as

$$x = A \cos \omega t. \quad (12)$$

Equation (12) describes the behaviour sketched graphically in Figure 3. The mass is initially displaced a distance $x = A$ and released at time $t = 0$. As time goes on, the mass oscillates from A to $-A$ and back to A again in the time it takes ωt to advance by 2π . This time is called T , the period of oscillation, so that $\omega T = 2\pi$, or $T = 2\pi/\omega$. The reciprocal of the period, or the frequency f , in oscillations per second, is given by $f = 1/T = \omega/2\pi$. The quantity ω is called the angular frequency and is expressed in radians per second.

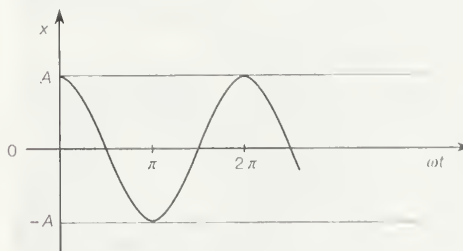


Figure 3: The function $x = A \cos \omega t$.

The choice of equation (12) as a possible kind of behaviour satisfying the differential equation (11) can be tested by substituting it into equation (11). The first derivative of x with respect to t is

$$\begin{aligned} \frac{dx}{dt} &= \frac{d}{dt}(A \cos \omega t) \\ &= -\omega A \sin \omega t. \end{aligned} \quad (13)$$

Differentiating a second time gives

$$\begin{aligned} \frac{d^2x}{dt^2} &= \frac{d}{dt} \left(\frac{dx}{dt} \right) & (14) \\ &= \frac{d}{dt} (-\omega A \sin \omega t) \\ &= -\omega^2 A \cos \omega t \\ &= -\omega^2 x. \end{aligned}$$

Equation (14) is the same as equation (11) if

$$\omega^2 = \frac{k}{m}. \quad (15)$$

Thus, subject to this condition, equation (12) is a correct solution to the differential equation. There are other possible correct guesses (e.g., $x = A \sin \omega t$) that differ from this one only in whether the mass is at rest or in motion at the instant $t = 0$.

The mass, as has been shown, oscillates from A to $-A$ and back again. The speed, given by dx/dt , equation (13), is zero at A and $-A$, but has its maximum magnitude, equal to ωA , when x is equal to zero. Physically, after the mass is displaced from equilibrium a distance A to the right, the restoring force F pushes the mass back toward its equilibrium position, causing it to accelerate to the left. When it reaches equilibrium, there is no force acting on it at that instant, but it is moving at speed ωA , and its inertia takes it past the equilibrium position. Before it is stopped it reaches position $-A$, and by this time there is a force acting on it again, pushing it back toward equilibrium.

The whole process, known as simple harmonic motion, repeats itself endlessly with a frequency given by equation (15). Equation (15) means that the stiffer the springs (i.e., the larger k), the higher the frequency (the faster the oscillations). Making the mass greater has exactly the opposite effect, slowing things down.

One of the most important features of harmonic motion is the fact that the frequency of the motion, ω (or f), depends only on the mass and the stiffness of the spring. It does not depend on the amplitude A of the motion. If the amplitude is increased, the mass moves faster, but the time required for a complete round trip remains the same. This fact has profound consequences, governing the nature of music and the principle of accurate timekeeping.

The potential energy of a harmonic oscillator, equal to the work an outside agent must do to push the mass from zero to x , is $U = \frac{1}{2}kx^2$. Thus, the total initial energy in the situation described above is $\frac{1}{2}kA^2$; and since the kinetic energy is always $\frac{1}{2}mv^2$, when the mass is at any point x in the oscillation,

$$\frac{1}{2}mv^2 + \frac{1}{2}kx^2 = \frac{1}{2}kA^2. \quad (16)$$

Equation (16) plays exactly the role for harmonic oscillators that equation (8) does for falling bodies.

It is quite generally true that harmonic oscillations result from disturbing any body or structure from a state of stable mechanical equilibrium. To understand this point, a brief discussion of stability is useful.

Consider a bowl with a marble resting inside, then consider a second, inverted bowl with a marble balanced on top. In both cases, the net force on the marble is zero. The marbles are thus in mechanical equilibrium. However, a small disturbance in the position of the marble balanced on top of the inverted bowl will cause it to roll away and not return. In such a case, the equilibrium is said to be unstable. Conversely, if the marble inside the first bowl is disturbed, gravity acts to push it back toward the bottom of the bowl. The marble inside the bowl (like the mass held by springs in Figure 2A) is an example of a body in stable equilibrium. If it is disturbed slightly, it executes harmonic oscillations around the bottom of the bowl rather than rolling away.

This argument may be generalized by a simple mathematical argument. Consider a body or structure in mechanical equilibrium, which, when disturbed by a small amount x , finds a force acting on it that is a function of x , $F(x)$. For small x , such a function may be written generally as a power series in x ; i.e.,

$$F(x) = F(0) + ax + bx^2 + \dots \quad (17)$$

where $F(0)$ is the value of $F(x)$ when $x = 0$, and a and b are constants, independent of x , determined by the nature of the system. The statement that the body is in mechanical equilibrium means that $F(0) = 0$, so that no force is acting on the body when it is undisturbed. Since x is small, x^2 is much smaller; thus the term bx^2 and all higher powers may be disregarded. This leaves $F(x) = ax$. Now, if a is positive, a disturbance produces a force in the same direction as the disturbance. This was the case when the marble was balanced on top of the inverted bowl. It describes unstable equilibrium. For the system to be stable, a must be negative. Thus, if $a = -k$, where k is some positive constant, equation (17) becomes $F(x) = -kx$, which is simply Hooke's law, equation (10). As has been described above, any system obeying Hooke's law is a harmonic oscillator.

The generality of this argument accounts for the fact that harmonic oscillators are abundantly observed in common experience. For example, any rigid structure will oscillate at many different harmonic frequencies corresponding to different possible distortions of its equilibrium shape. In addition, music may be produced either by disturbing the equilibrium of a stretched wire or fibre (as in the piano and violin), a stretched membrane (e.g., drums), or a rigid bar (the triangle and the xylophone) or by disturbing the density of an enclosed column of air (as in the trumpet and organ). While a fluid such as air is not rigid, its density is an example of a stable system that obeys Hooke's law and may therefore be set into harmonic oscillations.

All music would be quite different from what it is were it not for the general property of harmonic oscillators that the frequency is independent of the amplitude. Thus, instruments yield the same note (frequency) regardless of how loudly they are played (amplitude), and, equally important, the same note persists as the vibrations die away. This same property of harmonic oscillators is the underlying principle of all accurate timekeeping.

The first precise timekeeping mechanism, whose principles of motion were discovered by Galileo, was the simple pendulum (see below). The accuracy of modern timekeeping has been improved dramatically by the introduction of tiny quartz crystals, whose harmonic oscillations generate electrical signals that may be incorporated into miniaturized circuits in clocks and wristwatches. All harmonic oscillators are natural timekeeping devices because they oscillate at intrinsic natural frequencies independent of amplitude. A given number of complete cycles always corresponds to the same elapsed time. Quartz crystal oscillators make more accurate clocks than pendulums do principally because they oscillate many more times per second.

DAMPED AND FORCED OSCILLATIONS

The simple harmonic oscillations discussed above continue forever, at constant amplitude, oscillating as shown in Figure 3 between A and $-A$. Common experience indicates that real oscillators behave somewhat differently, however. Harmonic oscillations tend to die away as time goes on. This behaviour, called damping of the oscillations, is produced by forces such as friction and viscosity. These forces are known collectively as dissipative forces because they tend to dissipate the potential and kinetic energies of macroscopic bodies into the energy of the chaotic motion of atoms and molecules known as heat.

Friction and viscosity are complicated phenomena whose effects cannot be represented accurately by a general equation. However, for slowly moving bodies, the dissipative forces may be represented by

$$F_d = -\gamma v. \quad (18)$$

where v is the speed of the body and γ is a constant coefficient, independent of dynamic quantities such as speed or displacement. Equation (18) is most easily understood by an argument analogous to that applied to equation (17) above. F_d is written as a sum of powers of v , or $F_d(v) = F_d(0) + av + bv^2 + \dots$. When the body is at rest ($v = 0$), no dissipative force is expected because, if

Frequency of simple harmonic motion

Harmonic oscillators as time-keepers

Mechanical equilibrium

there were one, it might set the body into motion. Thus, $F_d(0) = 0$. The next term must be negative since dissipative forces always resist the motion. Thus, $a = -\gamma$ where γ is positive. Since v^2 has the same sign regardless of the direction of the motion, b must equal 0 lest it sometimes contribute a dissipative force in the same direction as the motion. The next term is proportional to v^3 , and it and all subsequent terms may be neglected if v is sufficiently small. So, as in equation (17) the power series is reduced to a single term, in this case $F_d = -\gamma v$.

To find the effect of a dissipative force on a harmonic oscillator, a new differential equation must be solved. The net force, or mass times acceleration, written as $m d^2x/dt^2$, is set equal to the sum of the Hooke's law force, $-kx$, and the dissipative force, $-\gamma v = -\gamma dx/dt$. Dividing by m yields

$$\frac{d^2x}{dt^2} = -\frac{k}{m}x - \frac{\gamma}{m} \frac{dx}{dt} \tag{19}$$

The general solution to equation (19) is given in the form $x = Ce^{-\gamma t/2m} \cos(\omega t + \theta_0)$, where C and θ_0 are arbitrary constants determined by the initial conditions. This motion, for the case in which $\theta_0 = 0$, is illustrated in Figure 4. As expected, the harmonic oscillations die out with time. The amplitude of the oscillations is bounded by an exponentially decreasing function of time (the dashed curves). The characteristic decay time (after which the oscillations are smaller by $1/e$, where e is the base of the natural logarithms $e = 2.718 \dots$) is equal to $2m/\gamma$. The frequency of the oscillations is given by

$$\omega^2 = \frac{k}{m} - \frac{\gamma^2}{4m^2} \tag{20}$$

Importantly, this frequency does not change as the oscillations decay.

From S.C. Frautschi et al., *The Mechanical Universe: Mechanics and Heat*, Advanced Edition (1986), Cambridge University Press

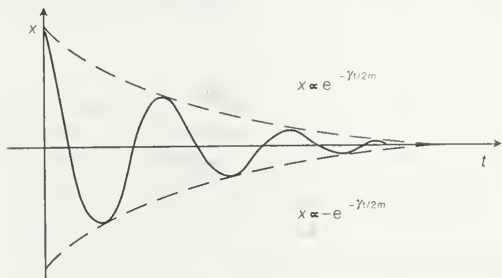


Figure 4: Damped oscillations.

Equation (20) shows that it is possible, by proper choice of γ , to turn a harmonic oscillator into a system that does not oscillate at all—that is, a system whose natural frequency is $\omega = 0$. Such a system is said to be critically damped. For example, the springs that suspend the body of an automobile cause it to be a natural harmonic oscillator. The shock absorbers of the auto are devices that seek to add just enough dissipative force to make the assembly critically damped. In this way, the passengers need not go through numerous oscillations after each bump in the road.

A simple disturbance can set a harmonic oscillator into motion. Repeated disturbances can increase the amplitude of the oscillations if they are applied in synchrony with the natural frequency. Even a very small disturbance, repeated periodically at just the right frequency, can cause a very large amplitude motion to build up. This phenomenon is known as resonance.

Periodically forced oscillations may be represented mathematically by adding a term of the form $a_0 \sin \omega t$ to the right-hand side of equation (19). This term describes a force applied at frequency ω , with amplitude ma_0 . The result of applying such a force is to create a kind of motion that does not need to decay with time, since the energy lost to dissipative processes is replaced, over the course of each cycle, by the driving force. The amplitude of the motion depends on how close the driving frequency ω is to the natural frequency ω_0 of the oscillator. Interestingly, even though dissipation is present, ω_0 is not given by

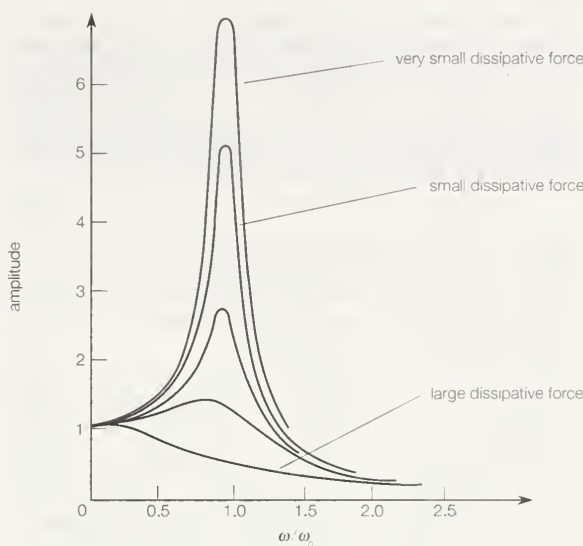


Figure 5: Resonance curves for different amounts of dissipative force acting on an oscillator.

From S.C. Frautschi et al., *The Mechanical Universe: Mechanics and Heat*, Advanced Edition (1986), Cambridge University Press

equation (20) but rather by equation (15): $\omega_0^2 = k/m$. The amplitude of the steady state motion (*i.e.*, long after the driving force has begun to be applied) is shown in Figure 5. The maximum amplitude occurs as expected at $\omega = \omega_0$. The height and width of the resonance curve are governed by the damping coefficient γ . If there were no damping, the maximum amplitude would be infinite. Because small disturbances at every possible frequency are always present in the natural world, every rigid structure would shake itself to pieces if not for the presence of internal damping.

Resonances are not uncommon in the world of familiar experience. For example, cars often rattle at certain engine speeds, and windows sometimes rattle when an airplane flies by. Resonance is particularly important in music. For example, the sound box of a violin does its job well if it has a natural frequency of oscillation that responds resonantly to each musical note. Very strong resonances to certain notes—called “wolf notes” by musicians—occur in cheap violins and are much to be avoided. Sometimes, a glass may be broken by a singer as a result of its resonant response to a particular musical note.

Motion of a particle in two or more dimensions

PROJECTILE MOTION

Galileo was quoted above pointing out with some detectable pride that none before him had realized that the curved path followed by a missile or projectile is a parabola. He had arrived at his conclusion by realizing that a body undergoing ballistic motion executes, quite independently, the motion of a freely falling body in the vertical direction and inertial motion in the horizontal direction. These considerations, and terms such as ballistic and projectile, apply to a body that, once launched, is acted upon by no force other than the Earth's gravity.

Projectile motion may be thought of as an example of motion in space—that is to say, of three-dimensional motion rather than motion along a line, or one-dimensional motion. In a suitably defined system of Cartesian coordinates, the position of the projectile at any instant may be specified by giving the values of its three coordinates, $x(t)$, $y(t)$, and $z(t)$. By generally accepted convention, $z(t)$ is used to describe the vertical direction. To a very good approximation, the motion is confined to a single vertical plane, so that for any single projectile it is possible to choose a coordinate system such that the motion is two-dimensional [say, $x(t)$ and $z(t)$] rather than three-dimensional [$x(t)$, $y(t)$, and $z(t)$]. It is assumed throughout this section that the range of the motion is sufficiently limited that the curvature of the Earth's surface may be ignored.

Consider a body whose vertical motion obeys equation (4), Galileo's law of falling bodies, which states $z = z_0 - 1/2gt^2$, while, at the same time, moving horizontally at

Resonance

Three-dimensional motion

a constant speed v_x in accordance with Galileo's law of inertia. The body's horizontal motion is thus described by $x(t) = v_x t$, which may be written in the form $t = x/v_x$. Using this result to eliminate t from equation (4) gives $z = z_0 - \frac{1}{2}g(1/v_x)^2 x^2$. This latter is the equation of the trajectory of a projectile in the z - x plane, fired horizontally from an initial height z_0 . It has the general form

$$z = a + bx^2, \tag{21}$$

where a and b are constants. Equation (21) may be recognized to describe a parabola (Figure 6A), just as Galileo claimed. The parabolic shape of the trajectory is preserved even if the motion has an initial component of velocity in the vertical direction (Figure 6B).

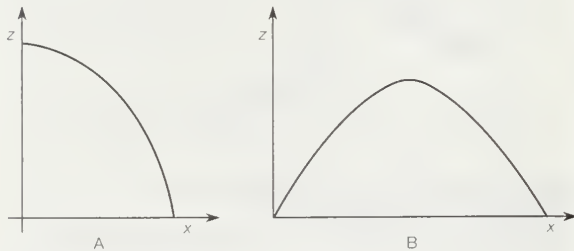


Figure 6: (A) The parabolic path of a projectile. (B) The parabolic path of a projectile with an initial upward component of velocity.

Energy is conserved in projectile motion. The potential energy $U(z)$ of the projectile is given by $U(z) = mgz$. The kinetic energy K is given by $K = \frac{1}{2}mv^2$, where v^2 is equal to the sum of the squares of the vertical and horizontal components of velocity, or $v^2 = v_x^2 + v_z^2$.

In all of this discussion, the effects of air resistance (to say nothing of wind and other more complicated phenomena) have been neglected. These effects are seldom actually negligible. They are most nearly so for bodies that are heavy and slow-moving. All of this discussion, therefore, is of great value for understanding the underlying principles of projectile motion but of little utility for predicting the actual trajectory of, say, a cannonball once fired or even a well-hit baseball.

Effects of air resistance

MOTION OF A PENDULUM

According to legend, Galileo discovered the principle of the pendulum while attending mass at the Duomo (cathedral) located in the Piazza del Duomo of Pisa, Italy. A lamp hung from the ceiling by a cable and, having just been lit, was swaying back and forth. Galileo realized that each complete cycle of the lamp took the same amount of time, compared to his own pulse, even though the amplitude of each swing was smaller than the last. As has already been shown, this property is common to all harmonic oscillators, and, indeed, Galileo's discovery led directly to the invention of the first accurate mechanical clocks. Galileo was also able to show that the period of oscillation of a simple pendulum is proportional to the square root of its length and does not depend on its mass.

Galileo's discovery of pendulum motion

A simple pendulum is sketched in Figure 7. A bob of mass M is suspended by a massless cable or bar of length L from a point about which it pivots freely. The angle between the cable and the vertical is called θ . The force

From R.P. Olenick, T.M. Apostol, and D.L. Goodstein, *The Mechanical Universe: Introduction to Mechanics and Heat* (1985) Cambridge University Press

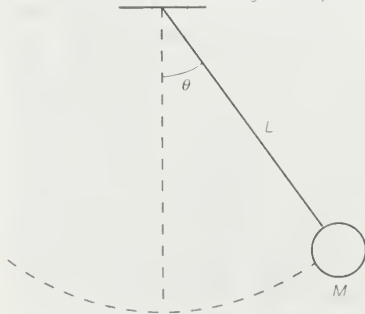


Figure 7: A simple pendulum (see text).

of gravity acting on the mass M , always equal to $-Mg$ in the vertical direction, is a vector that may be resolved into two components, one that acts ineffectually along the cable and another, perpendicular to the cable, that tends to restore the bob to its equilibrium position directly below the point of suspension. This latter component is given by

$$F = -Mg \sin \theta. \tag{22}$$

The bob is constrained by the cable to swing through an arc that is actually a segment of a circle of radius L . If the cable is displaced through an angle θ , the bob moves a distance $L\theta$ along its arc (θ must be expressed in radians for this form to be correct). Thus, Newton's second law may be written

$$F = Ma = M \frac{d^2(L\theta)}{dt^2}. \tag{23}$$

Equating equation (22) to equation (23), one sees immediately that the mass M will drop out of the resulting equation. The simple pendulum is an example of a falling body, and its dynamics do not depend on its mass for exactly the same reason that the acceleration of a falling body does not depend on its mass: both the force of gravity and the inertia of the body are proportional to the same mass, and the effects cancel one another. The equation that results (after extracting the constant L from the derivative and dividing both sides by L) is

$$\frac{d^2\theta}{dt^2} = -\frac{g}{L} \sin \theta. \tag{24}$$

If the angle θ is sufficiently small, equation (24) may be rewritten in a form that is both more familiar and more amenable to solution. Figure 8 shows a segment of a circle of radius L . A radius vector at angle θ , as shown, locates a point on the circle displaced a distance $L\theta$ along the arc.

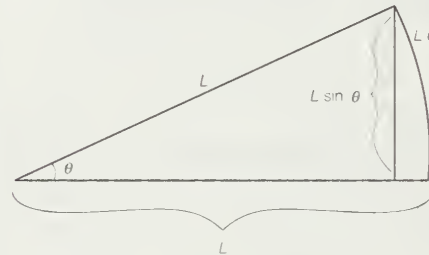


Figure 8: A segment of a circle of radius L (see text).

It is clear from the geometry that $L \sin \theta$ and $L\theta$ are very nearly equal for small θ . It follows then that $\sin \theta$ and θ are also very nearly equal for small θ . Thus, if the analysis is restricted to small angles, then $\sin \theta$ may be replaced by θ in equation (24) to obtain

$$\frac{d^2\theta}{dt^2} = -\frac{g}{L} \theta. \tag{25}$$

Equation (25) should be compared with equation (11): $d^2x/dt^2 = -(k/m)x$. In the first case, the dynamic variable (meaning the quantity that changes with time) is θ , in the second case it is x . In both cases, the second derivative of the dynamic variable with respect to time is equal to the variable itself multiplied by a negative constant. The equations are therefore mathematically identical and have the same solution—i.e., equation (12), or $\theta = A \cos \omega t$. In the case of the pendulum, the frequency of the oscillations is given by the constant in equation (25), or $\omega^2 = g/L$. The period of oscillation, $T = 2\pi/\omega$, is therefore

$$T = 2\pi \sqrt{\frac{L}{g}}.$$

Period of oscillation

Just as Galileo concluded, the period is independent of the mass and proportional to the square root of the length.

As with most problems in physics, this discussion of the pendulum has involved a number of simplifications and approximations. Most obviously, $\sin \theta$ was replaced by θ to obtain equation (25). This approximation is surprisingly accurate. For example, at a not-very-small angle of 17.2° , corresponding to 0.300 radian, $\sin \theta$ is equal to 0.296, an

error of less than 2 percent. For smaller angles, of course, the error is appreciably smaller.

The problem was also treated as if all the mass of the pendulum were concentrated at a point at the end of the cable. This approximation assumes that the mass of the bob at the end of the cable is much larger than that of the cable and that the physical size of the bob is small compared with the length of the cable. When these approximations are not sufficient, one must take into account the way in which mass is distributed in the cable and bob. This is called the physical pendulum, as opposed to the idealized model of the simple pendulum. Significantly, the period of a physical pendulum does not depend on its total mass either.

The effects of friction, air resistance, and the like have also been ignored. These dissipative forces have the same effects on the pendulum as they do on any other kind of harmonic oscillator, as discussed above. They cause the amplitude of a freely swinging pendulum to grow smaller on successive swings. Conversely, in order to keep a pendulum clock going, a mechanism is needed to restore the energy lost to dissipative forces.

CIRCULAR MOTION

Consider a particle moving along the perimeter of a circle at a uniform rate, such that it makes one complete revolution every hour. To describe the motion mathematically, a vector is constructed from the centre of the circle to the particle. The vector then makes one complete revolution every hour. In other words, the vector behaves exactly like the large hand on a wristwatch, an arrow of fixed length that makes one complete revolution every hour. The motion of the point of the vector is an example of uniform circular motion, and the period T of the motion is equal to one hour ($T = 1$ h). The arrow sweeps out an angle of 2π radians (one complete circle) per hour. This rate is called the angular frequency and is written $\omega = 2\pi \text{ h}^{-1}$. Quite generally, for uniform circular motion at any rate,

$$T = \frac{2\pi}{\omega} \tag{26}$$

These definitions and relations are the same as they are for harmonic motion, discussed above.

Consider a coordinate system, as shown in Figure 9A, with the circle centred at the origin. At any instant of time, the position of the particle may be specified by giving the radius r of the circle and the angle θ between the position vector and the x -axis. Although r is constant, θ increases uniformly with time t , such that $\theta = \omega t$, or $d\theta/dt = \omega$, where ω is the angular frequency in equation (26). Contrary to the case of the wristwatch, however, ω is positive by convention when the rotation is in the counterclockwise sense. The vector r has x and y components given by

$$x = r \cos \theta = r \cos \omega t, \tag{27}$$

$$y = r \sin \theta = r \sin \omega t. \tag{28}$$

One meaning of equations (27) and (28) is that, when a particle undergoes uniform circular motion, its x and y components each undergo simple harmonic motion. They are, however, not in phase with one another: at the instant when x has its maximum amplitude (say, at $\theta = 0$), y has zero amplitude, and vice versa.

In a short time, Δt , the particle moves $r\Delta\theta$ along the circumference of the circle, as shown in Figure 9B. The average speed of the particle is thus given by

$$\bar{v} = r \frac{\Delta\theta}{\Delta t} \tag{29}$$

The average velocity of the particle is a vector given by

$$\bar{\mathbf{v}} = \frac{\mathbf{r}(t + \Delta t) - \mathbf{r}(t)}{\Delta t} \tag{30}$$

This operation of vector subtraction is indicated in Figure 9B. It yields a vector that is nearly perpendicular to $r(t)$ and $r(t + \Delta t)$. Indeed, the instantaneous velocity, found by allowing Δt to shrink to zero, is a vector \mathbf{v} that is perpendicular to r at every instant and whose magnitude is

$$|\mathbf{v}| = r \frac{d\theta}{dt} = r\omega. \tag{31}$$

The relationship between r and \mathbf{v} is shown in Figure 9C. It means that the particle's instantaneous velocity is always tangent to the circle.

Notice that, just as the position vector r may be described in terms of the components x and y given by equations (27) and (28), the velocity vector \mathbf{v} may be described in terms of its projections on the x and y axes, given by

$$v_x = \frac{dx}{dt} = -r\omega \sin \omega t, \tag{32}$$

$$v_y = \frac{dy}{dt} = r\omega \cos \omega t. \tag{33}$$

Imagine a new coordinate system, in which a vector of length ωr extends from the origin and points at all times in the same direction as \mathbf{v} . This construction is shown in Figure 9D. Each time the particle sweeps out a complete circle, this vector also sweeps out a complete circle. In fact, its point is executing uniform circular motion at the same angular frequency as the particle itself. Because vectors have magnitude and direction, but not position in space, the vector that has been constructed is the velocity \mathbf{v} . The velocity of the particle is itself undergoing uniform circular motion at angular frequency ω .

Although the speed of the particle is constant, the parti-

From (A) R.P. Olenick, T.M. Apostol, and D.L. Goodstein, *The Mechanical Universe Introduction to Mechanics and Heat* (1985) Cambridge University Press

Angular frequency

Average particle velocity

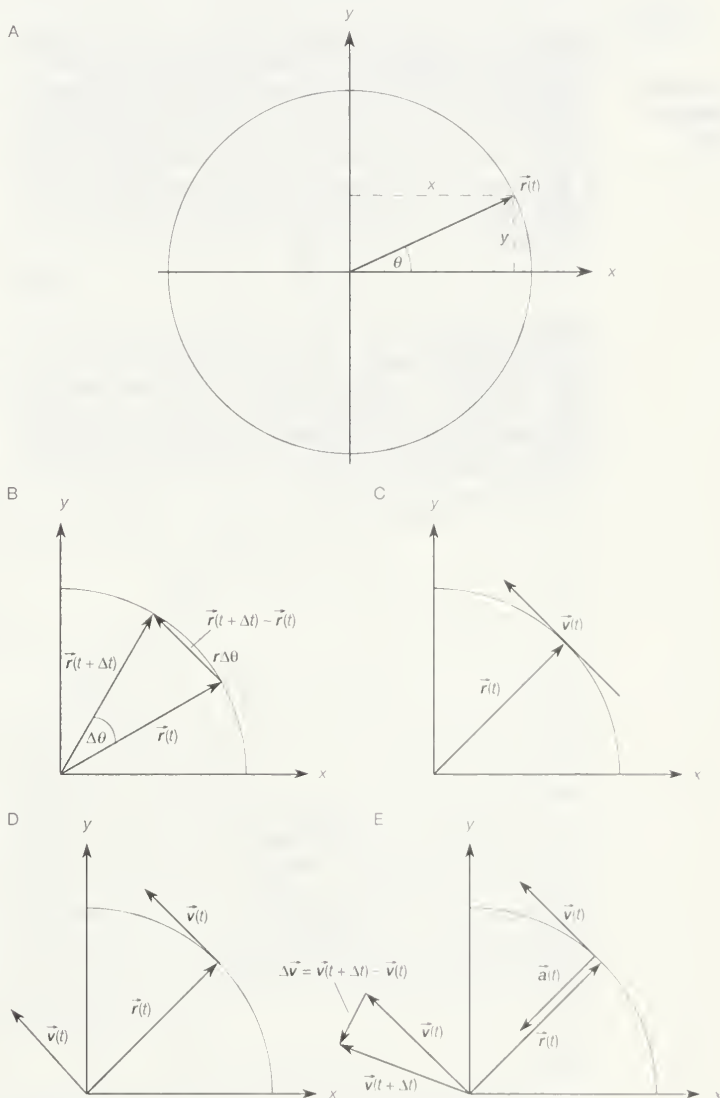


Figure 9: (A) A coordinate system to describe uniform circular motion. (B) The distance traveled in time Δt by a particle undergoing uniform circular motion. (C) The instantaneous velocity of the particle. (D) The velocity vector \mathbf{v} undergoes uniform circular motion at the same angular frequency as the particle. (E) The acceleration vector of the particle. (See text.)

cle is nevertheless accelerated, because its velocity is constantly changing direction. The acceleration \mathbf{a} is given by

$$\mathbf{a} = \frac{d\mathbf{v}}{dt}. \quad (34)$$

Since \mathbf{v} is a vector of length $r\omega$ undergoing uniform circular motion, equations (29) and (30) may be repeated, as illustrated in Figure 9E, giving

$$\bar{a} = r\omega \frac{\Delta\theta}{\Delta t} \quad (35)$$

$$\bar{\mathbf{a}} = \frac{\mathbf{v}(t + \Delta t) - \mathbf{v}(t)}{\Delta t}. \quad (36)$$

Thus, one may conclude that the instantaneous acceleration is always perpendicular to \mathbf{v} and its magnitude is

$$|\mathbf{a}| = r\omega \frac{d\theta}{dt} = r\omega^2. \quad (37)$$

Since \mathbf{v} is perpendicular to \mathbf{r} , and \mathbf{a} is perpendicular to \mathbf{v} , the vector \mathbf{a} is rotated 180° with respect to \mathbf{r} . In other words, the acceleration is parallel to \mathbf{r} but in the opposite direction. The same conclusion may be reached by realizing that \mathbf{a} has x and y components given by

$$a_x = \frac{dv_x}{dt} = -r\omega^2 \cos \omega t, \quad (38)$$

$$a_y = \frac{dv_y}{dt} = -r\omega^2 \sin \omega t, \quad (39)$$

similar to equations (32) and (33). When equations (38) and (39) are compared with equations (27) and (28) for x and y , it is clear that the components of \mathbf{a} are just those of \mathbf{r} multiplied by $-\omega^2$, so that $\mathbf{a} = -\omega^2\mathbf{r}$. This acceleration is called the centripetal acceleration, meaning that it is inward, pointing along the radius vector toward the centre of the circle. It is sometimes useful to express the centripetal acceleration in terms of the speed v . Using $v = \omega r$, one can write

$$a = -\frac{v^2}{r}. \quad (40)$$

CIRCULAR ORBITS

The detailed behaviour of real orbits is the concern of celestial mechanics (see below). This section treats only the idealized, uniform circular orbit of a planet such as the Earth about a central body such as the Sun. In fact, the Earth's orbit about the Sun is not quite exactly uniformly

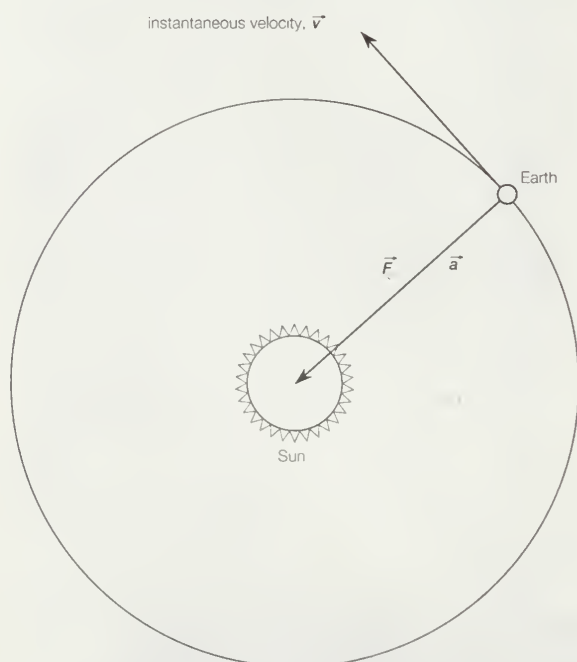


Figure 10: The gravitational force F_G exerted by the Sun on the Earth produces the centripetal acceleration \mathbf{a}_c of the Earth's orbital motion.

circular, but it is a close enough approximation for the purposes of this discussion.

A body in uniform circular motion undergoes at all times a centripetal acceleration given by equation (40). According to Newton's second law, a force is required to produce this acceleration. In the case of an orbiting planet, the force is gravity. The situation is illustrated in Figure 10. The gravitational attraction of the Sun is an inward (centripetal) force acting on the Earth. This force produces the centripetal acceleration of the orbital motion.

Before these ideas are expressed quantitatively, an understanding of why a force is needed to maintain a body in an orbit of constant speed is useful. The reason is that, at each instant, the velocity of the planet is tangent to the orbit. In the absence of gravity, the planet would obey the law of inertia (Newton's first law) and fly off in a straight line in the direction of the velocity at constant speed. The force of gravity serves to overcome the inertial tendency of the planet, thereby keeping it in orbit.

The gravitational force between two bodies such as the Sun and the Earth is given by

$$F = -G \frac{M_S M_E}{r^2}, \quad (41)$$

where M_S and M_E are the masses of the Sun and the Earth, respectively, r is the distance between their centres, and G is a universal constant equal to $6.672 \times 10^{-11} \text{ Nm}^2/\text{kg}^2$ (Newton metres squared per kilogram squared). The force acts along the direction connecting the two bodies (*i.e.*, along the radius vector of the uniform circular motion), and the minus sign signifies that the force is attractive, acting to pull the Earth toward the Sun.

To an observer on the surface of the Earth, the planet appears to be at rest at (approximately) a constant distance from the Sun. It would appear to the observer, therefore, that any force (such as the Sun's gravity) acting on the Earth must be balanced by an equal and opposite force that keeps the Earth in equilibrium. In other words, if gravity is trying to pull the Earth into the Sun, some opposing force must be present to prevent that from happening. In reality, no such force exists. The Earth is in freely accelerated motion caused by an unbalanced force. The apparent force, known in mechanics as a pseudoforce, is due to the fact that the observer is actually in accelerated motion. In the case of orbital motion, the outward pseudoforce that balances gravity is called the centrifugal force.

For a uniform circular orbit, gravity produces an inward acceleration given by equation (40), $a = -v^2/r$. The pseudoforce f needed to balance this acceleration is just equal to the mass of the Earth times an equal and opposite acceleration, or $f = M_E v^2/r$. The earthbound observer then believes that there is no net force acting on the planet—*i.e.*, that $F + f = 0$, where F is the force of gravity given by equation (41). Combining these equations yields a relation between the speed v of a planet and its distance r from the Sun:

$$v^2 = G \frac{M_S}{r}. \quad (42)$$

It should be noted that the speed does not depend on the mass of the planet. This occurs for exactly the same reason that all bodies fall toward Earth with the same acceleration and that the period of a pendulum is independent of its mass. An orbiting planet is in fact a freely falling body.

Equation (42) is a special case (for circular orbits) of Kepler's third law, which will be encountered below in *Celestial mechanics*. Using the fact that $v = 2\pi r/T$, where $2\pi r$ is the circumference of the orbit and T is the time to make a complete orbit (*i.e.*, T is one year in the life of the planet), it is easy to show that $T^2 = (4\pi^2/GM_S)r^3$. This relation also may be applied to satellites in circular orbit around the Earth (in which case, M_E must be substituted for M_S) or in orbit around any other central body.

ANGULAR MOMENTUM AND TORQUE

A particle of mass m and velocity \mathbf{v} has linear momentum $\mathbf{p} = m\mathbf{v}$. The particle may also have angular momentum \mathbf{L} with respect to a given point in space. If \mathbf{r} is the vector from the point to the particle, then

$$\mathbf{L} = \mathbf{r} \times \mathbf{p}. \quad (43)$$

Centripetal
acceleration

Pseudo-
force

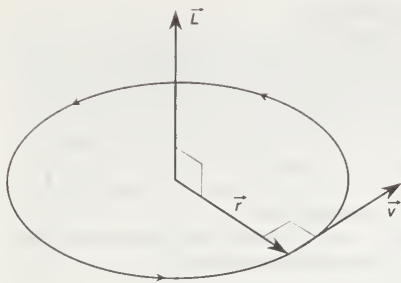


Figure 11: The angular momentum L of a particle traveling in a circular orbit.

From R.P. Olenick, T.M. Apostol, and D.L. Goodstein, *The Mechanical Universe. Introduction to Mechanics and Heat* (1985). Cambridge University Press.

Notice that angular momentum is always a vector perpendicular to the plane defined by the vectors r and p (or v). For example, if the particle (or a planet) is in a circular orbit, its angular momentum with respect to the centre of the circle is perpendicular to the plane of the orbit and in the direction given by the vector cross product right-hand rule, as shown in Figure 11. Moreover, since in the case of a circular orbit, r is perpendicular to p (or v), the magnitude of L is simply

$$L = rp = mvr. \quad (44)$$

The significance of angular momentum arises from its derivative with respect to time,

$$\frac{dL}{dt} = \frac{d}{dt}(r \times p) = m \frac{d}{dt}(r \times v), \quad (45)$$

where p has been replaced by mv and the constant m has been factored out. Using the product rule of differential calculus,

$$\frac{d}{dt}(r \times v) = \frac{dr}{dt} \times v + r \times \frac{dv}{dt}. \quad (46)$$

In the first term on the right-hand side of equation (46), dr/dt is simply the velocity v , leaving $v \times v$. Since the cross product of any vector with itself is always zero, that term drops out, leaving

$$\frac{d}{dt}(r \times v) = r \times \frac{dv}{dt}. \quad (47)$$

Here, dv/dt is the acceleration a of the particle. Thus, if equation (47) is multiplied by m , the left-hand side becomes dL/dt , as in equation (45), and the right-hand side may be written $r \times ma$. Since, according to Newton's second law, ma is equal to F , the net force acting on the particle, the result is

$$\frac{dL}{dt} = r \times F. \quad (48)$$

Equation (48) means that any change in the angular momentum of a particle must be produced by a force that is not acting along the same direction as r . One particularly important application is the solar system. Each planet is held in its orbit by its gravitational attraction to the Sun, a force that acts along the vector from the Sun to the planet. Thus the force of gravity cannot change the angular momentum of any planet with respect to the Sun. Therefore, each planet has constant angular momentum with respect to the Sun. This conclusion is correct even though the real orbits of the planets are not circles but ellipses.

The quantity $r \times F$ is called the torque τ . Torque may be thought of as a kind of twisting force, the kind needed to tighten a bolt or to set a body into rotation. Using this definition, equation (48) may be rewritten

$$\tau = r \times F = \frac{dL}{dt}. \quad (49)$$

Equation (49) means that if there is no torque acting on a particle, its angular momentum is constant, or conserved. Suppose, however, that some agent applies a force F_a to the particle resulting in a torque equal to $r \times F_a$. According to Newton's third law, the particle must apply a force $-F_a$ to the agent. Thus there is a torque equal to $-r \times F_a$ acting on the agent. The torque on the particle

causes its angular momentum to change at a rate given by $dL/dt = r \times F_a$. However, the angular momentum L_a of the agent is changing at the rate $dL_a/dt = -r \times F_a$. Therefore, $dL/dt + dL_a/dt = 0$, meaning that the total angular momentum of particle plus agent is constant, or conserved. This principle may be generalized to include all interactions between bodies of any kind, acting by way of forces of any kind. Total angular momentum is always conserved. The law of conservation of angular momentum is one of the most important principles in all of physics.

Motion of a group of particles

CENTRE OF MASS

The word particle has been used in this article to signify an object whose entire mass is concentrated at a point in space. In the real world, however, there are no particles of this kind. All real bodies have sizes and shapes. Furthermore, as Newton believed and is now known, all bodies are in fact compounded of smaller bodies called atoms. Therefore, the science of mechanics must deal not only with particles but also with more complex bodies that may be thought of as collections of particles.

To take a specific example, the orbit of a planet around the Sun was discussed earlier as if the planet and the Sun were each concentrated at a point in space. In reality, of course, each is a substantial body. However, because each is nearly spherical in shape, it turns out to be permissible, for the purposes of this problem, to treat each body as if its mass were concentrated at its centre. This is an example of an idea that is often useful in discussing bodies of all kinds: the centre of mass. The centre of mass of a uniform sphere is located at the centre of the sphere. For many purposes (such as the one cited above) the sphere may be treated as if all its mass were concentrated at its centre of mass.

To extend the idea further, consider the Earth and the Sun not as two separate bodies but as a single system of two bodies interacting with one another by means of the force of gravity. In the previous discussion of circular orbits, the Sun was assumed to be at rest at the centre of the orbit, but, according to Newton's third law, it must actually be accelerated by a force due to the Earth that is equal and opposite to the force that the Sun exerts on the Earth. In other words, considering only the Sun and Earth (ignoring, for example, all the other planets), if M_S and M_E are, respectively, the masses of the Sun and the Earth, and if a_S and a_E are their respective accelerations, then combining Newton's second and third laws results in the equation $M_S a_S = -M_E a_E$. Writing each a as dv/dt , this equation is easily manipulated to give

$$\frac{d}{dt}(M_S v_S + M_E v_E) = 0, \quad (50)$$

or

$$M_S v_S + M_E v_E = \text{constant}. \quad (51)$$

This remarkable result means that, as the Earth orbits the Sun and the Sun moves in response to the Earth's gravitational attraction, the entire two-body system has constant linear momentum, moving in a straight line at constant speed. Without any loss of generality, one can imagine observing the system from a frame of reference moving along with that same speed and direction. This is sometimes called the centre-of-mass frame. In this frame, the momentum of the two-body system—*i.e.*, the constant in equation (51)—is equal to zero. Writing each of the v 's as the corresponding dr/dt , equation (51) may be expressed in the form

$$\frac{d}{dt}(M_S r_S + M_E r_E) = 0. \quad (52)$$

Thus, $M_S r_S$ and $M_E r_E$ are two vectors whose vector sum does not change with time. The sum is defined to be the constant vector MR , where M is the total mass of the system and equals $M_S + M_E$. Thus,

$$MR = M_S r_S + M_E r_E. \quad (53)$$

Centre of mass of a sphere

This procedure defines a constant vector \mathbf{R} , from any arbitrarily chosen point in space. The relation between vectors \mathbf{R} , \mathbf{r}_S , and \mathbf{r}_E is shown in Figure 12. The fact that \mathbf{R} is constant (although \mathbf{r}_S and \mathbf{r}_E are not constant) means that, rather than the Earth orbiting the Sun, the Earth and Sun are both orbiting an imaginary point fixed in space. This point is known as the centre of mass of the two-body system.

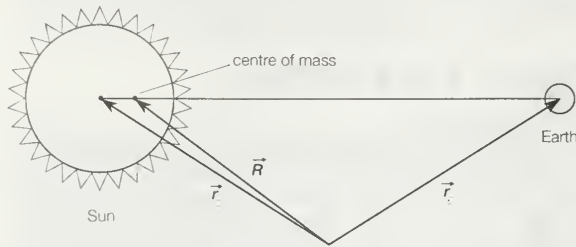


Figure 12: The centre of mass of the two-body Earth-Sun system.

Knowing the masses of the two bodies ($M_S = 1.99 \times 10^{30}$ kilograms, $M_E = 5.98 \times 10^{24}$ kilograms), it is easy to find the position of the centre of mass. The origin of the coordinate system may be chosen to be located at the centre of mass merely by defining $\mathbf{R} = 0$. Then $r_S = (M_E/M_S) r_E \approx 450$ kilometres, when r_E is rounded to 1.5×10^8 km. A few hundred kilometres is so small compared to r_E that, for all practical purposes, no appreciable error occurs when r_S is ignored and the Sun is assumed to be stationary at the centre of the orbit.

Centre of mass for an N -body system

With this example as a guide, it is now possible to define the centre of mass of any collection of bodies. Assume that there are N bodies altogether, each labeled with numbers ranging from 1 to N , and that the vector from an arbitrary origin to the i th body—where i is some number between 1 and N —is \mathbf{r}_i , as shown in Figure 13. Let the mass of the i th body be m_i . Then the total mass of the N -body system is

$$m = \sum_{i=1}^N m_i \tag{54}$$

and the centre of mass of the system is found at the end of a vector \mathbf{R} given by

$$m\mathbf{R} = \sum_{i=1}^N m_i \mathbf{r}_i \tag{55}$$

as illustrated in Figure 13. This definition applies regardless of whether the N bodies making up the system are the stars in a galaxy, the atoms in a rigid body, larger and arbitrarily chosen segments of a rigid body, or any other system of masses. According to equation (55), the vector to the centre of mass of any system is a kind of weighted average of the vectors to all the components of the system.

As will be demonstrated in the sections that follow, the statics and dynamics of many complicated bodies or systems may often be understood by simply applying Newton's laws as if the system's mass were concentrated at the centre of mass.

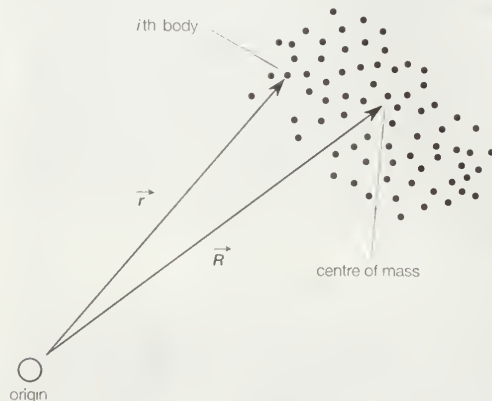


Figure 13: The centre of mass of an N -body system (see text).

CONSERVATION OF MOMENTUM

Newton's second law, in its most general form, says that the rate of a change of a particle's momentum \mathbf{p} is given by the force acting on the particle; i.e., $\mathbf{F} = d\mathbf{p}/dt$. If there is no force acting on the particle, then, since $d\mathbf{p}/dt = 0$, \mathbf{p} must be constant, or conserved. This observation is merely a restatement of Newton's first law, the principle of inertia: if there is no force acting on a body, it moves at constant speed in a straight line.

Now suppose that an external agent applies a force \mathbf{F}_a to the particle so that \mathbf{p} changes according to

$$\frac{d\mathbf{p}}{dt} = \mathbf{F}_a \tag{56}$$

According to Newton's third law, the particle must apply an equal and opposite force $-\mathbf{F}_a$ to the external agent. The momentum \mathbf{p}_a of the external agent therefore changes according to

$$\frac{d\mathbf{p}_a}{dt} = -\mathbf{F}_a \tag{57}$$

Adding together equations (56) and (57) results in the equation

$$\frac{d}{dt}(\mathbf{p} + \mathbf{p}_a) = 0 \tag{58}$$

The force applied by the external agent changes the momentum of the particle, but at the same time the momentum of the external agent must also change in such a way that the total momentum of both together is constant, or conserved. This idea may be generalized to give the law of conservation of momentum: in all the interactions between all the bodies in the universe, total momentum is always conserved.

Conservation of momentum

It is useful in this light to examine the behaviour of a complicated system of many parts. The centre of mass of the system may be found using equation (55). Differentiating with respect to time gives

$$m\mathbf{v} = \sum_{i=1}^N m_i \mathbf{v}_i \tag{59}$$

where $\mathbf{v} = d\mathbf{R}/dt$ and $\mathbf{v}_i = d\mathbf{r}_i/dt$. Note that $m_i \mathbf{v}_i$ is the momentum of the i th part of the system, and $m\mathbf{v}$ is the momentum that the system would have if all its mass (i.e., m) were concentrated at its centre of mass, the point whose velocity is \mathbf{v} . Thus, the momentum associated with the centre of mass is the sum of the momenta of the parts.

Suppose now that there is no external agent applying a force to the entire system. Then the only forces acting on the system are those exerted by the parts on one another. These forces may accelerate the individual parts. Differentiating equation (59) with respect to time gives

$$m \frac{d\mathbf{v}}{dt} = \sum_{i=1}^N m_i \frac{d\mathbf{v}_i}{dt} = \sum_{i=1}^N \mathbf{F}_i \tag{60}$$

where \mathbf{F}_i is the net force, or the sum of the forces, exerted by all the other parts of the body on the i th part. \mathbf{F}_i is defined mathematically by the equation

$$\mathbf{F}_i = \sum_{j=1}^N \mathbf{F}_{ij} \tag{61}$$

where \mathbf{F}_{ij} represents the force on body i due to body j (the force on body i due to itself, \mathbf{F}_{ii} , is zero). The motion of the centre of mass is then given by the complicated-looking formula

$$m \frac{d\mathbf{v}}{dt} = \sum_{i=1}^N \left(\sum_{j=1}^N \mathbf{F}_{ij} \right) \tag{62}$$

This complicated formula may be greatly simplified, however, by noting that Newton's third law requires that for every force \mathbf{F}_{ij} exerted by the j th body on the i th body, there is an equal and opposite force $-\mathbf{F}_{ij}$ exerted by the i th body on the j th body. In other words, every term in the double sum has an equal and opposite term. The double summation on the right-hand side of equation (61) always adds up to zero. This result is true regardless of the complexity of the system, the nature of the forces acting between the parts, or the motions of the parts. In short, in the absence of external forces acting

on the system as a whole, $mdv/dt = 0$, which means that the momentum of the centre of mass of the system is always conserved. Having determined that momentum is conserved whether or not there is an external force acting, one may conclude that the total momentum of the universe is always conserved.

COLLISIONS

A collision is an encounter between two bodies that alters at least one of their courses. Altering the course of a body requires that a force be applied to it. Thus, each body exerts a force on the other. These forces of interaction may operate at some distance, as do the gravitational and electromagnetic forces, or the bodies may appear to make physical contact. However, even apparent contact between two bodies is only a macroscopic manifestation of microscopic forces that act between atoms some distance apart. There is no fundamental distinction between physical contact and interaction at a distance.

The importance of understanding the mechanics of collisions is obvious to anyone who has ever driven an automobile. In modern physics, however, collisions are important for a different reason. The current understanding of the subatomic particles of which atoms are composed is derived entirely from studying the results of collisions among them. Thus, in modern physics, the description of collisions is a significant part of the understanding of matter. These descriptions are quantum mechanical rather than classical, but they are nevertheless closely based on principles that arise out of classical mechanics.

It is possible in principle to predict the result of a collision using Newton's second law directly. Suppose that two bodies are going to collide and that F , the force of interaction between them, is known to be a function of r , the distance between them. Then, if it is known that, say, one particle has incident momentum p , the problem is solved if the final momentum $p + \Delta p$ can be determined. Inverting Newton's second law, $F = dp/dt$, the change in momentum is given by

$$\Delta p = \int_{-\infty}^{\infty} F dt. \tag{63}$$

This integral is known as the impulse imparted to the particle. In order to perform the integral, it is necessary to know r at all times so that F may be known at all times. More realistically, Δp is the sum of a series of small steps, such that

$$\delta p = F \delta t, \tag{64}$$

where F depends on the instantaneous distance between the particles. Because $p = mv = mdr/dt$, the change in r in this step is

$$\delta r = \frac{p}{m} \delta t. \tag{65}$$

At the next step, there is a new distance, $r + \delta r$, giving a new value of the force in equation (64) and a new momentum, $p + \delta p$, in equation (65). This method of analyzing collisions is used in numerical calculations on digital computers.

To predict the result of a collision analytically (rather than numerically) it is often most useful to apply conservation laws. In any collision (as in any other phenomenon), energy, momentum, and angular momentum are always conserved. Judicious application of these laws may be extremely useful because they do not depend in any way on the detailed nature of the interaction (*i.e.*, the force as a function of distance).

This point can be illustrated by the following example. A collision is to take place between two bodies of the same mass m . One of the bodies is initially at rest (its momentum is zero). The other has initial momentum p_0 . After the collision, the body previously at rest has momentum p_1 , and the body initially in motion has momentum p_2 . Since momentum is conserved, the total momentum after the collision, $p_1 + p_2$, must be equal to the total momentum before the collision, p_0 ; that is,

$$p_0 = p_1 + p_2. \tag{66}$$

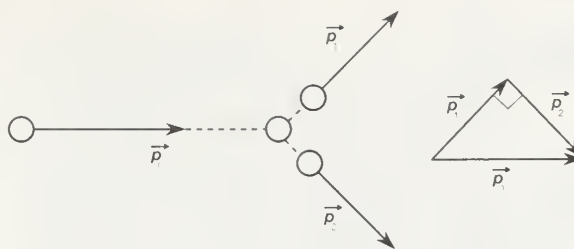


Figure 14: Collision between two particles of equal mass.

Equation (66) is the equation of a vector triangle, as shown in Figure 14. However, p_1 and p_2 are not determined by this condition; they are only constrained by it.

Although energy is always conserved, the kinetic energy of the incident body is not always converted entirely into the kinetic energy of the two bodies after the collision. For example, if the bodies are microscopic (say, two identical atoms), the collision may cause one or both to be excited into a state of higher internal energy than it started with. Such an event would leave correspondingly less kinetic energy for the outgoing atoms. In fact, it is precisely by studying the trajectories of outgoing projectiles in collisions like these that physicists are able to determine the possible excited states of microscopic particles.

In a collision between macroscopic objects, some of the kinetic energy is always converted to heat. Heat is the energy of random vibrations of the atoms and molecules that constitute the bodies. However, if the amount of heat is negligible compared to the initial kinetic energy, it may be ignored. Such a collision is said to be elastic.

Suppose the collision described above between two bodies, each of mass m , is between billiard balls, and suppose it is elastic (a reasonably good approximation of real billiard balls). The kinetic energy of the incident ball is then equal to the sum of the kinetic energies of the outgoing balls. According to equation (3), the kinetic energy of a moving object is given by $K = \frac{1}{2}mv^2$, where v is the speed of the ball (technically, the energy associated with the fact that the ball is rolling as well as translating is ignored here; see below *Rotation about a moving axis*). Equation (3) may be written in a particularly useful form by recognizing that since $p = mv$

$$K = \frac{1}{2}mv^2 = \frac{p^2}{2m}. \tag{67}$$

Then the conservation of kinetic energy may be written

$$\frac{p_0^2}{2m} = \frac{p_1^2}{2m} + \frac{p_2^2}{2m}, \tag{68}$$

or, canceling the factors $2m$,

$$p_0^2 = p_1^2 + p_2^2. \tag{69}$$

Comparing this result with equation (66) shows that the vector triangle is pythagorean; p_1 and p_2 are perpendicular. This result is well known to all experienced pool players. Notice that it was possible to arrive at this result without any knowledge of the forces that act when billiard balls collide.

RELATIVE MOTION

A collision between two bodies can always be described in a frame of reference in which the total momentum is zero. This is the centre-of-mass (or centre-of-momentum) frame mentioned earlier. Then, for example, in the collision between two bodies of the same mass discussed above, the two bodies always have equal and opposite velocities, as shown in Figure 15. It should be noted that, in this frame of reference, the outgoing momenta are antiparallel and not perpendicular.

Any collection of bodies may similarly be described in a frame of reference in which the total momentum is zero. This frame is simply the one in which the centre of mass is at rest. This fact is easily seen by differentiating equation (55) with respect to time, giving

$$m \frac{dR}{dt} = \sum_{i=1}^N m_i \frac{dr_i}{dt}. \tag{70}$$

Significance of understanding collision mechanics

Elastic collisions

Impulse

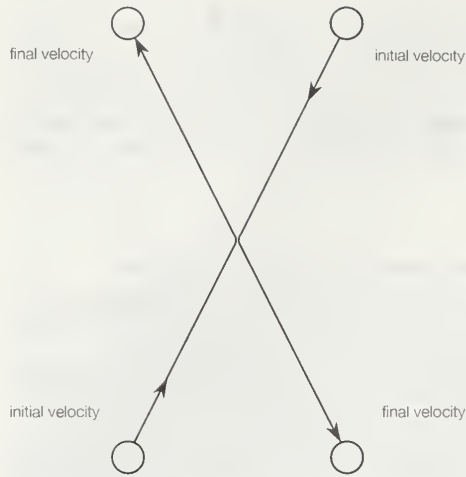


Figure 15: Collision between two particles of equal mass as seen from the centre-of-mass frame of reference.

The right-hand side is the sum of the momenta of all the bodies. It is equal to zero if the velocity of the centre of mass, $d\mathbf{R}/dt$, is equal to zero.

Principle of Galilean relativity

If Newton's second law is correct in any frame of reference, it will also appear to be correct to an observer moving with any constant velocity with respect to that frame. This principle, called the principle of Galilean relativity, is true because, to the moving observer, the same constant velocity seems to have been added to the velocity of every particle in the system. This change does not affect the accelerations of the particles (since the added velocity is constant, not accelerated) and therefore does not change the apparent force (mass times acceleration) acting on each particle. That is why it is permissible to describe a problem from the centre-of-momentum frame (provided that the centre of mass is not accelerated) or from any other frame moving at constant velocity with respect to it.

If this principle is strictly correct, the fundamental forces of physics should not contain any particular speed. This must be true because the speed of any object will be different to observers in different but equally good frames of reference, but the force should always be the same. It turns out, according to the theory of James Clerk Maxwell, that there is an intrinsic speed in the force laws of electricity and magnetism: the speed of light appears in the forces between electric charges and between magnetic poles. This discrepancy was ultimately resolved by Albert Einstein's special theory of relativity. According to the special theory of relativity, Newtonian mechanics breaks down when the relative speed between particles approaches the speed of light (see below *Relativistic mechanics*).

COUPLED OSCILLATORS

In the section on simple harmonic oscillators, the motion of a single particle held in place by springs was considered. In this section, the motion of a group of particles bound by springs to one another is discussed. The solutions of this seemingly academic problem have far-reaching implications in many fields of physics. For example, a system of particles held together by springs turns out to be a useful model of the behaviour of atoms mutually bound in a crystalline solid.

To begin with a simple case, consider two particles in a line, as shown in Figure 16. Each particle has mass m , each spring has spring constant k , and motion is restricted

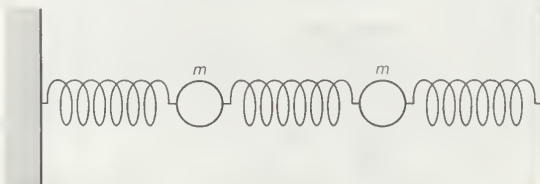


Figure 16: Coupled oscillators (see text).

to the horizontal, or x , direction. Even this elementary system is capable of surprising behaviour, however. For instance, if one particle is held in place while the other is displaced, and then both are released, the displaced particle immediately begins to execute simple harmonic motion. This motion, by stretching the spring between the particles, starts to excite the second particle into motion. Gradually the energy of motion passes from the first particle to the second until a point is reached at which the first particle is at rest and only the second is oscillating. Then the process starts all over again, the energy passing in the opposite direction.

To analyze the possible motions of the system, one writes equations similar to equation (11), giving the acceleration of each particle owing to the forces acting on it. There is one equation for each particle (two equations in this case). The force on each particle depends not only on its displacement from its equilibrium position but also on its distance from the other particle, since the spring between them stretches or compresses according to that distance. For this reason the motions are coupled, the solution of each equation (the motion of each particle) depending on the solution of the other (the motion of the other).

Analyzing the system yields the fact that there are two special states of motion in which both particles are always in oscillation with the same frequency. In one state, the two particles oscillate in opposite directions with equal and opposite displacements from equilibrium at all times. In the other state, both particles move together, so that the spring between them is never stretched or compressed. The first of these motions has higher frequency than the second because the centre spring contributes an increase in the restoring force.

These two collective motions, at different, definite frequencies, are known as the normal modes of the system.

Normal modes

If a third particle is inserted into the system together with another spring, there will be three equations to solve, and the result will be three normal modes. A large number N of particles in a line will have N normal modes. Each normal mode has a definite frequency at which all the particles oscillate. In the highest frequency mode each particle moves in the direction opposite to both of its neighbours. In the lowest frequency mode, neighbours move almost together, barely disturbing the springs between them. Starting from one end, the amplitude of the motion gradually builds up, each particle moving a bit more than the one before, reaching a maximum at the centre, and then decreasing again. A plot of the amplitudes, shown in Figure 17, basically describes one-half of a sine wave from one end of the system to the other. The next mode is $3/2$ of a sine wave, and so on to the highest frequency mode, which may be visualized as $(2N - 1)/2$ sine waves. If the vibrations were up and down rather than side to side, these modes would be identical to the fundamental and harmonic vibrations excited by plucking a guitar string.

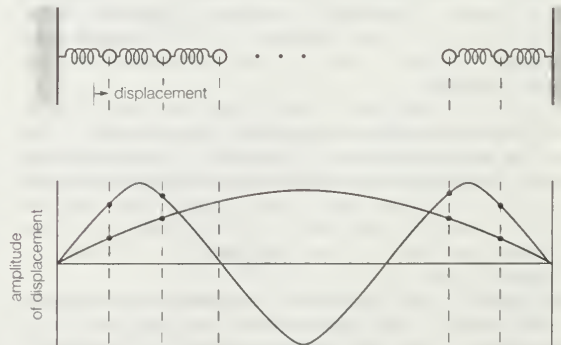


Figure 17: Normal modes (see text).

The atoms of a crystal are held in place by mutual forces of interaction that oppose any disturbance from equilibrium positions, just as the spring forces in the example above. For small displacements of the atoms, they behave mathematically just like spring forces—*i.e.*, they obey Hooke's law, equation (10). Each atom is free to move in three dimensions rather than one, however;

Excitation of modes of a crystal

therefore each atom added to a crystal adds three normal modes. In a typical crystal at ordinary temperature, all these modes are always excited by random thermal energy. The lower-frequency, longer-wavelength modes may also be excited mechanically. These are called sound waves.

Rigid bodies

STATICS

Statics is the study of bodies and structures that are in equilibrium. For a body to be in equilibrium, there must be no net force acting on it. In addition, there must be no net torque acting on it. Figure 18A shows a body in equilibrium under the action of equal and opposite forces. Figure 18B shows a body acted on by equal and opposite forces that produce a net torque, tending to start it rotating. It is therefore not in equilibrium.

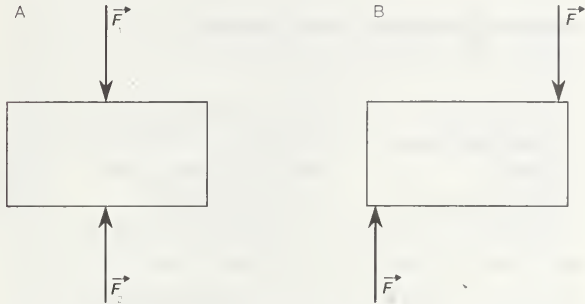


Figure 18: (A) A body in equilibrium under equal and opposite forces. (B) A body not in equilibrium under equal and opposite forces.

When a body has a net force and a net torque acting on it owing to a combination of forces, all the forces acting on the body may be replaced by a single (imaginary) force called the resultant, which acts at a single point on the body, producing the same net force and the same net torque. The body can be brought into equilibrium by applying to it a real force at the same point, equal and opposite to the resultant. This force is called the equilibrant. An example is shown in Figure 19.

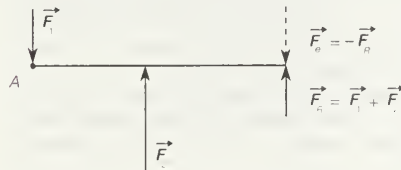


Figure 19: The resultant force (F_r) produces the same net force and the same net torque about point A as $F_1 + F_2$; the body can be brought into equilibrium by applying the equilibrant force F_e .

The torque on a body due to a given force depends on the reference point chosen, since the torque τ by definition equals $r \times F$, where r is a vector from some chosen reference point to the point of application of the force. Thus, for a body to be at equilibrium, not only must the net force on it be equal to zero but the net torque with respect to any point must also be zero. Fortunately, it is easily shown for a rigid body that, if the net force is zero and the net torque is zero with respect to any one point, then the net torque is also zero with respect to any other point in the frame of reference.

Rigid bodies

A body is formally regarded as rigid if the distance between any set of two points in it is always constant. In reality no body is perfectly rigid. When equal and opposite forces are applied to a body, it is always deformed slightly. The body's own tendency to restore the deformation has the effect of applying counterforces to whatever is applying the forces, thus obeying Newton's third law. Calling a body rigid means that the changes in the dimensions of the body are small enough to be neglected, even though the force produced by the deformation may not be neglected.

Equal and opposite forces acting on a rigid body may act so as to compress the body (Figure 20A) or to stretch it



Figure 20: (A) Compression produced by equal and opposite forces. (B) Tension produced by equal and opposite forces.

(Figure 20B). The bodies are then said to be under compression or under tension, respectively. Strings, chains, and cables are rigid under tension but may collapse under compression. On the other hand, certain building materials, such as brick and mortar, stone, or concrete, tend to be strong under compression but very weak under tension.

Application of statics

The most important application of statics is to study the stability of structures, such as edifices and bridges. In these cases, gravity applies a force to each component of the structure as well as to any bodies the structure may need to support. The force of gravity acts on each bit of mass of which each component is made, but for each rigid component it may be thought of as acting at a single point, the centre of gravity, which is in these cases the same as the centre of mass.

To give a simple but important example of the application of statics, consider the two situations shown in Figure 21. In each case, a mass m is supported by two symmetric members, each making an angle θ with respect to the horizontal. In Figure 21A the members are under tension; in Figure 21B they are under compression. In either case, the force acting along each of the members is shown to be

$$F = \frac{mg}{2 \sin \theta} \quad (71)$$

The force in either case thus becomes intolerably large if the angle θ is allowed to be very small. In other words, the mass cannot be hung from a perfectly horizontal member.

The ancient Greeks built magnificent stone temples; however, the horizontal stone slabs that constituted the roofs of the temples could not support even their own weight over more than a very small span. For this reason, one characteristic that identifies a Greek temple is the many closely spaced pillars needed to hold up the flat roof. The problem posed by equation (71) was solved by the ancient Romans, who incorporated into their architecture the arch, a structure that supports its weight by compression, corresponding to Figure 21B.

A suspension bridge illustrates the use of tension. The weight of the span and any traffic on it is supported by cables, which are placed under tension by the weight. Corresponding to Figure 21A, the cables are not stretched to be horizontal, but rather they are always hung so as to have substantial curvature.

It should be mentioned in passing that equilibrium under static forces is not sufficient to guarantee the stability of a structure. It must also be stable against perturbations such as the additional forces that might be imposed, for example, by winds or by earthquakes. Analysis of the stability of structures under such perturbations is an important part of the job of an engineer or architect.

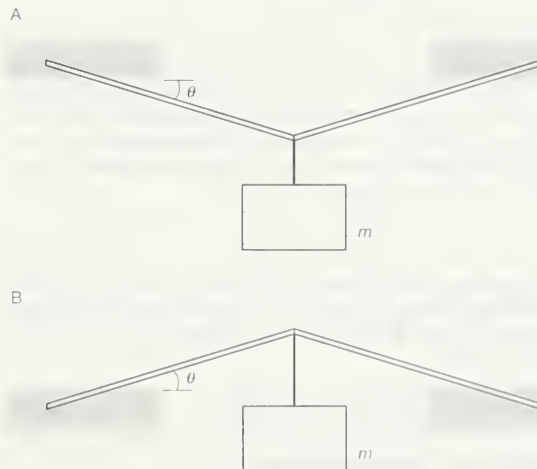


Figure 21: (A) A body supported by two rigid members under tension. (B) A body supported by two rigid members under compression.

ROTATION ABOUT A FIXED AXIS

Consider a rigid body that is free to rotate about an axis fixed in space. Because of the body's inertia, it resists being set into rotational motion, and equally important, once rotating, it resists being brought to rest. Exactly how that inertial resistance depends on the mass and geometry of the body is discussed here.

Factors affecting inertial resistance

Take the axis of rotation to be the z -axis. A vector in the x - y plane from the axis to a bit of mass fixed in the body makes an angle θ with respect to the x -axis. If the body is rotating, θ changes with time, and the body's angular frequency is

$$\omega = \frac{d\theta}{dt}; \tag{72}$$

ω is also known as the angular velocity. If ω is changing in time, there is also an angular acceleration α , such that

$$\alpha = \frac{d\omega}{dt}. \tag{73}$$

Because linear momentum p is related to linear speed v by $p = mv$, where m is the mass, and because force F is related to acceleration a by $F = ma$, it is reasonable to assume that there exists a quantity I that expresses the rotational inertia of the rigid body in analogy to the way m expresses the inertial resistance to changes in linear motion. One would expect to find that the angular momentum is given by

$$L = I\omega \tag{74}$$

and that the torque (twisting force) is given by

$$\tau = I\alpha. \tag{75}$$

One can imagine dividing the rigid body into bits of mass labeled m_1, m_2, m_3 , and so on. Let the bit of mass at the tip of the vector be called m_i , as indicated in Figure 22.

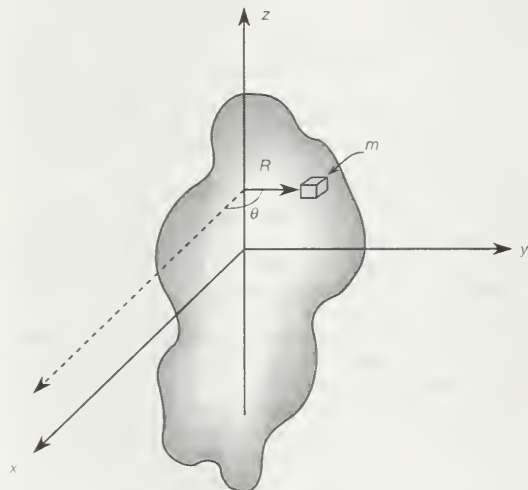


Figure 22: Rotation around a fixed axis.

If the length of the vector from the axis to this bit of mass is R_i , then m_i 's linear velocity v_i equals ωR_i (see equation [31]), and its angular momentum L_i equals $m_i v_i R_i$ (see equation [44]), or $m_i R_i^2 \omega$. The angular momentum of the rigid body is found by summing all the contributions from all the bits of mass labeled $i = 1, 2, 3 \dots$:

$$L = \left(\sum_i m_i R_i^2 \right) \omega. \tag{76}$$

In a rigid body, the quantity in parentheses in equation (76) is always constant (each bit of mass m_i always remains the same distance R_i from the axis). Thus if the motion is accelerated, then

$$\frac{dL}{dt} = \left(\sum_i m_i R_i^2 \right) \frac{d\omega}{dt}. \tag{77}$$

Recalling that $\tau = dL/dt$, one may write

$$\tau = \left(\sum_i m_i R_i^2 \right) \alpha. \tag{78}$$

(These equations may be written in scalar form, since L and τ are always directed along the axis of rotation in this discussion.) Comparing equations (76) and (78) with (74) and (75), one finds that

$$I = \sum_i m_i R_i^2. \tag{79} \text{ Moment of inertia}$$

The quantity I is called the moment of inertia.

According to equation (79), the effect of a bit of mass on the moment of inertia depends on its distance from the axis. Because of the factor R_i^2 , mass far from the axis makes a bigger contribution than mass close to the axis. It is important to note that R_i is the distance from the axis, not from a point. Thus, if x_i and y_i are the x and y coordinates of the mass m_i , then $R_i^2 = x_i^2 + y_i^2$, regardless of the value of the z coordinate. The moments of inertia of some simple uniform bodies are given in the table.

Moments of Inertia for Uniform Bodies		
body	axis	I
Thin rod (length L)	perpendicular axis through centre	$\frac{1}{12}ML^2$
Thin ring (radius R)	perpendicular axis through centre	MR^2
Solid circular cylinder	axis of cylinder	$\frac{1}{2}MR^2$
Thin disk	transverse axis through centre	$\frac{1}{4}MR^2$
Solid sphere	any axis through centre	$\frac{2}{5}MR^2$
Thin spherical shell	any axis through centre	$\frac{2}{3}MR^2$
Rectangular plate (length a , height b)	axis through centre perpendicular to the plate	$\frac{1}{12}M(a^2 + b^2)$

The moment of inertia of any body depends on the axis of rotation. Depending on the symmetry of the body, there may be as many as three different moments of inertia about mutually perpendicular axes passing through the centre of mass. If the axis does not pass through the centre of mass, the moment of inertia may be related to that about a parallel axis that does so. Let I_c be the moment of inertia about the parallel axis through the centre of mass, r the distance between the two axes, and M the total mass of the body. Then

$$I = I_c + Mr^2. \tag{80}$$

In other words, the moment of inertia about an axis that does not pass through the centre of mass is equal to the moment of inertia for rotation about an axis through the centre of mass (I_c) plus a contribution that acts as if the mass were concentrated at the centre of mass, which then rotates about the axis of rotation.

The dynamics of rigid bodies rotating about fixed axes may be summarized in three equations. The angular momentum is $L = I\omega$, the torque is $\tau = I\alpha$, and the kinetic energy is $K = \frac{1}{2}I\omega^2$.

ROTATION ABOUT A MOVING AXIS

The general motion of a rigid body tumbling through space may be described as a combination of translation of the body's centre of mass and rotation about an axis through the centre of mass. The linear momentum of the body of mass M is given by

$$\mathbf{p} = M\mathbf{v}_c \tag{81}$$

where \mathbf{v}_c is the velocity of the centre of mass. Any change in the momentum is governed by Newton's second law, which states that

$$\mathbf{F} = \frac{d\mathbf{p}}{dt}, \tag{82}$$

where \mathbf{F} is the net force acting on the body. The angular momentum of the body with respect to any reference point may be written as

$$\mathbf{L} = \mathbf{L}_c + \mathbf{r} \times \mathbf{p}, \tag{83}$$

where \mathbf{L}_c is the angular momentum of rotation about an axis through the centre of mass, \mathbf{r} is a vector from the reference point to the centre of mass, and $\mathbf{r} \times \mathbf{p}$ is therefore the angular momentum associated with motion of the centre of mass, acting as if all the body's mass were concentrated at that point. The quantity \mathbf{L}_c in equation (83) is sometimes called the body's spin, and $\mathbf{r} \times \mathbf{p}$ is called the orbital angular momentum. Any change in the angular momentum of the body is given by the torque equation.

Combined translation and rotation

$$\tau = \frac{dL}{dt}. \quad (84)$$

An example of a body that undergoes both translational and rotational motion is the Earth, which rotates about an axis through its centre once per day while executing an orbit around the Sun once per year. Because the Sun exerts no torque on the Earth with respect to its own centre, the orbital angular momentum of the Earth is constant in time. However, the Sun does exert a small torque on the Earth with respect to the planet's centre, owing to the fact that the Earth is not perfectly spherical. The result is a slow shifting of the Earth's axis of rotation, known as the precession of the equinoxes (see below).

The kinetic energy of a body that is both translating and rotating is given by

$$K = \frac{1}{2}Mv_c^2 + \frac{1}{2}I\omega^2, \quad (85)$$

where I is the moment of inertia and ω is the angular velocity of rotation about the axis through the centre of mass.

A common example of combined rotation and translation is rolling motion, as exhibited by a billiard ball rolling on a table, or a ball or cylinder rolling down an inclined plane. Consider the latter example, illustrated in Figure 23. Motion is impelled by the force of gravity, which may be resolved into two components, F_N , which is normal to the plane, and F_p , which is parallel to it. In addition to gravity, friction plays an essential role. The force of friction, written as f , acts parallel to the plane, in opposition to the direction of motion, at the point of contact between the plane and the rolling body. If f is very small, the body will slide without rolling. If f is very large, it will prevent motion from occurring. The magnitude of f depends on the smoothness and composition of the body and the plane, and it is proportional to F_N , the normal component of the force.

Frictional force

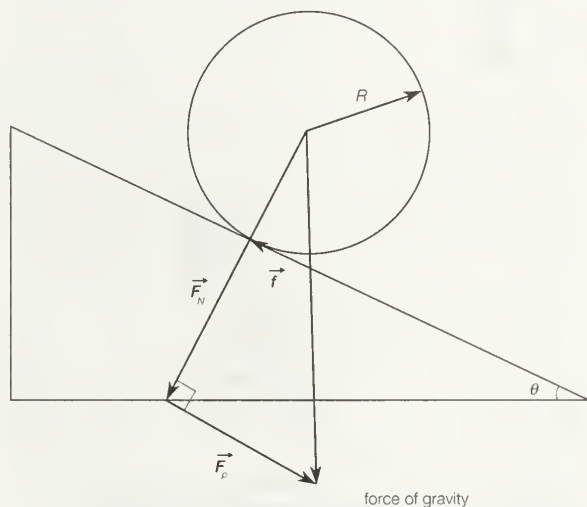


Figure 23: Rolling motion (see text).

Consider a case in which f is just large enough to cause the body (sphere or cylinder) to roll without slipping. The motion may be analyzed from the point of view of an axis passing through the point of contact between the rolling body and the plane. Remarkably, the point of contact may always be regarded to be instantaneously at rest. To understand why, suppose that the rolling body has radius R and angular velocity ω about its centre-of-mass axis. Then, with respect to its own axis, each point on the circular cross section in Figure 23 moves with instantaneous tangential linear speed $v_c = R\omega$. In particular, the point of contact is moving backward with this speed relative to the centre of mass. But with respect to the inclined plane, the centre of mass is moving forward with exactly this same speed. The net effect of the two equal and opposite speeds is that the point of contact is always instantaneously at rest. Therefore, although friction acts at that point, no work is done by friction, so mechanical energy (potential plus kinetic) may be regarded as conserved.

With respect to the axis through the point of contact, the torque is equal to RF_p , giving rise to an angular acceleration α given by $I_p\alpha = RF_p$, where I_p is the moment of inertia about the point-of-contact axis and can be determined by applying equation (80) relating moments of inertia about parallel axes ($I_p = I + MR^2$). Thus,

$$\alpha = \frac{RF_p}{I + MR^2}. \quad (86)$$

From this result, the motion of the body is easily obtained using the fact that the velocity of the centre of mass is $v_c = R\omega$ and hence the linear acceleration of the centre of mass is $a_c = R\alpha$.

Notice that, although without friction no angular acceleration would occur, the force of friction does not affect the magnitude of α . Because friction does no work, this same result may be obtained by applying energy conservation. The situation also may be analyzed entirely from the point of view of the centre of mass. In that case, the torque is $-fR$, but f also provides a linear force on the body. The f may then be eliminated by using Newton's second law and the fact that the torque equals the moment of inertia times the angular acceleration, once again leading to the same result.

One more interesting fact is hidden in the form of equation (86). The parallel component of the force of gravity is given by

$$F_p = Mg \sin \theta, \quad (87)$$

where θ is the angle of inclination of the plane. The moment of inertia about the centre of mass of any body of mass M may be written

$$I = Mk^2, \quad (88)$$

where k is a distance called the radius of gyration. Comparison to equation (79) shows that k is a measure of how far from the centre of mass the mass of the body is concentrated. Using equations (87) and (88) in equation (86), one finds that

Radius of gyration

$$\alpha = \frac{Rg \sin \theta}{k^2 + R^2}. \quad (89)$$

Thus, the angular acceleration of a body rolling down a plane does not depend on its total mass, although it does depend on its shape and distribution of mass. The same may be said of a_c , the linear acceleration of the centre of mass. The acceleration of a rolling ball, like the acceleration of a freely falling object, is independent of its mass. This observation helps to explain why Galileo was able to discover many of the basic laws of dynamics in gravity by studying the behaviour of balls rolling down inclined planes.

MOTION IN A ROTATING FRAME

Centrifugal force. According to the principle of Galilean relativity, if Newton's laws are true in any reference frame, they are also true in any other frame moving at constant velocity with respect to the first one. Conversely, they do not appear to be true in any frame accelerated with respect to the first. Instead, in an accelerated frame, objects appear to have forces acting on them that are not in fact present. These are called pseudoforces, as described above. Since rotational motion is always accelerated motion, pseudoforces may always be observed in rotating frames of reference.

As one example, a frame of reference in which the Earth is at rest must rotate once per year about the Sun. In this reference frame, the gravitational force attracting the Earth toward the Sun appears to be balanced by an equal and opposite outward force that keeps the Earth in stationary equilibrium. This outward pseudoforce, discussed above, is the centrifugal force.

The rotation of the Earth about its own axis also causes pseudoforces for observers at rest on the Earth's surface. There is a centrifugal force, but it is much smaller than the force of gravity. Its effect is that, at the Equator, where it is largest, the gravitational acceleration g is about 0.5 percent smaller than at the poles, where there is no centrifugal force. This same centrifugal force is responsible

for the fact that the Earth is slightly nonspherical, bulging just a bit at the Equator.

The tides

Pseudoforces can have real consequences. The oceanic tides on Earth, for example, are a consequence of centrifugal forces in the Earth-Moon and Earth-Sun systems. The Moon appears to be orbiting the Earth, but in reality both the Moon and the Earth orbit their common centre of mass. The centre of mass of the Earth-Moon system is located inside the Earth nearly three-fourths of the distance from the centre to the surface, or roughly 4,700 kilometres from the centre of the Earth. The Earth rotates about this point approximately once a month. The gravitational attraction of the Moon and the centrifugal force of this rotation are exactly balanced at the centre of the Earth. At the surface of the Earth closest to the Moon, the Moon's gravity is stronger than the centrifugal force. The ocean's waters, which are free to move in response to this unbalanced force, tend to build up a small bulge at that point. On the surface of the Earth exactly opposite the Moon, the centrifugal force is stronger than the Moon's gravity, and a small bulge of water tends to build up there as well. The water is correspondingly depleted at the points 90° on either side of these. Each day the Earth rotates beneath these bulges and troughs, which remain stationary with respect to the Earth-Moon system. The result is two high tides and two low tides every day every place on Earth. The Sun has a similar effect, but of only about half the size; it increases or decreases the size of the tides depending on its relative alignment with the Earth and Moon.

Coriolis force. The Coriolis force is a pseudoforce that operates in all rotating frames. One way to envision it is to imagine a rotating platform (such as a merry-go-round or a phonograph turntable) with a perfectly smooth surface and a smooth block sliding inertially across it. The block, having no (real) forces acting on it, moves in a straight line at constant speed in inertial space. However, the platform rotates under it, so that to an observer on the platform, the block appears to follow a curved trajectory, bending in the opposite direction to the motion of the platform. Since the motion is curved, and hence accelerated, there appears, to the observer, to be a force operating. That pseudoforce is called the Coriolis force.

The Coriolis force also may be observed on the surface of the Earth. For example, many science museums have a pendulum, called a Foucault pendulum, suspended from a long cable with markers to show that its plane of motion rotates slowly. The rotation of the plane of motion is caused by the Coriolis force. The effect is most easily imagined by picturing the pendulum swinging directly above the North Pole. The plane of its motion remains stationary in inertial space, while the Earth rotates once a day beneath it.

Effect on projectile motion

At lower latitudes, the effect is a bit more subtle, but it is still present. Imagine that, somewhere in the Northern Hemisphere, a projectile is fired due south. As viewed from inertial space, the projectile initially has an eastward component of velocity as well as a southward component because the gun that fired it, which is stationary on the surface of the Earth, was moving eastward with the Earth's rotation at the instant it was fired. However, since it was fired to the south, it lands at a slightly lower latitude, closer to the Equator. As one moves south, toward the Equator, the tangential speed of the Earth's surface due to its rotation increases because the surface is farther from the axis of rotation. Thus, although the projectile has an eastward component of velocity (in inertial space), it lands at a place where the surface of the Earth has a larger eastward component of velocity. Thus, to the observer on Earth, the projectile seems to curve slightly to the west. That westward curve is attributed to the Coriolis force. If the projectile were fired to the north, it would seem to curve eastward.

The same analysis applied to a Foucault pendulum explains why its plane of motion tends to rotate in the clockwise direction anywhere in the Northern Hemisphere and in the counterclockwise direction in the Southern Hemisphere. Storms, known as cyclones, tend to rotate in the opposite direction in each hemisphere, also due to the

Coriolis force. Air moves in all directions toward a low-pressure centre. In the Northern Hemisphere, air moving up from the south is deflected eastward, while air moving down from the north is deflected westward. This effect tends to give cyclones a counterclockwise circulation in the Northern Hemisphere. In the Southern Hemisphere, cyclones tend to circulate in the clockwise direction.

SPINNING TOPS AND GYROSCOPES

Figure 24A shows a wheel that is weighted in its rim to maximize its moment of inertia I and that is spinning with angular frequency ω on a horizontal axle supported at both ends. As shown, it has an angular momentum L along the x direction equal to $I\omega$. Now suppose the support at point P is removed, leaving the axle supported only at one end. Gravity, acting on the mass of the wheel as if it were concentrated at the centre of mass, applies a downward force on the wheel. The wheel, however, does not fall. Instead, the axle remains (nearly) horizontal but rotates in the counterclockwise direction as seen from above (Figure 24B). This motion is called gyroscopic precession.

From (B) R.P. Olenick, T.M. Apostol, and D.L. Goodstein, *The Mechanical Universe: Introduction to Mechanics and Heat* (1985), Cambridge University Press

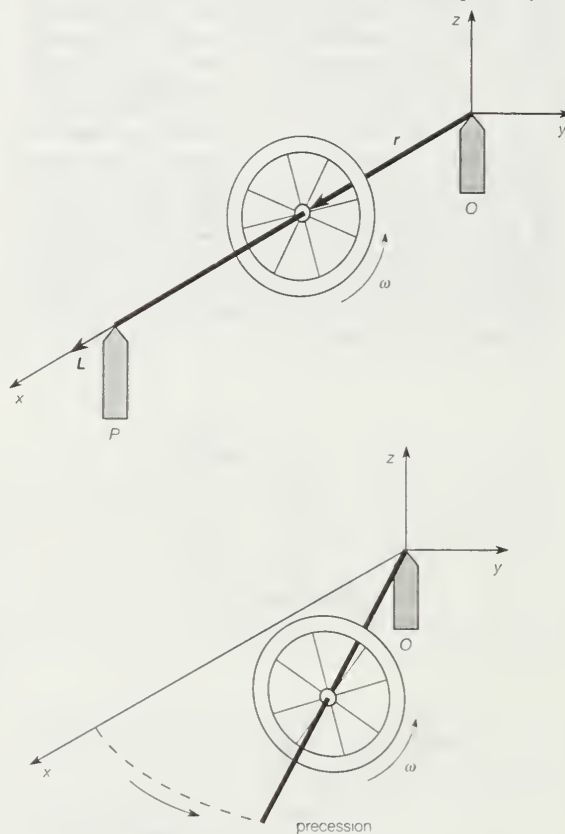


Figure 24: (A) A spinning wheel supported at both ends of an axle. (B) A spinning wheel supported at one end of an axle and exhibiting precession. (See text.)

Horizontal precession occurs in this case because the gravitational force results in a torque with respect to the point of suspension, such that $\tau = r \times F$ and is directed, initially, in the positive y direction. The torque causes the angular momentum L to move toward that direction according to $\tau = dL/dt$. Because τ is perpendicular to L , it does not change the magnitude of the angular momentum, only its direction. As precession proceeds, the torque remains horizontal, and the angular momentum vector, continually redirected by the torque, executes uniform circular motion in the horizontal plane at a frequency Ω , the frequency of precession.

Precession

In reality, the motion is a bit more complicated than uniform precession in the horizontal plane. When the support at P is released, the centre of mass of the wheel initially drops slightly below the horizontal plane. This drop reduces the gravitational potential energy of the system, releasing kinetic energy for the orbital motion of the

centre of mass as it precesses. It also provides a small component of L in the negative z direction, which balances the angular momentum in the positive z direction that results from the orbital motion of the centre of mass. There can be no net angular momentum in the vertical direction because there is no component of torque in that direction.

One more complication: the initial drop of the centre of mass carries it too far for a stable plane of precession, and it tends to bounce back up after overshooting. This produces an up-and-down oscillation during precession, called nutation ("nodding"). In most cases, nutation is quickly damped by friction in the bearings, leaving uniform precession.

A spinning top undergoes all the motions described above. If it is initially set spinning with a vertical axis, there will be virtually no torque, and conservation of angular momentum will keep the axis vertical for a long time. Eventually, however, friction at the point of contact will require the centre of mass to lower itself, which can only happen if the axis tilts. The spinning will also slow down, making the tilting process easier. Once the top tilts, gravity produces a horizontal torque that leads to precession of the spin axis. The subsequent motion depends on whether the point of contact is fixed or free to slip on the horizontal plane. Vast tomes have been written on the motions of tops.

A gyroscope is a device that is designed to resist changes in the direction of its axis of spin. That purpose is generally accomplished by maximizing its moment of inertia about the spin axis and by spinning it at the maximum practical frequency. Each of these considerations has the effect of maximizing the magnitude of the angular momentum, thus requiring a larger torque to change its direction. It is quite generally true that the torque τ , the angular momentum L , and the precession frequency Ω (defined as a vector along the precession axis in the direction given by the right-hand rule) are related by

$$\tau = \Omega \times L. \tag{90}$$

Equation (90), illustrated in Figure 25, is called the gyroscope equation.

Uses of gyroscopes

Gyroscopes are used for a variety of purposes, including navigation. Use of gyroscopes for this purpose is called inertial guidance. The gyroscope is suspended as nearly as possible at its centre of mass, so that gravity does not apply a torque that causes it to precess. The gyroscope tends therefore to point in a constant direction in space, allowing the orientation of the vehicle to be accurately maintained.

One further application of the gyroscope principle may be seen in the precession of the equinoxes. The Earth is a kind of gyroscope, spinning on its axis once each day. The Sun would apply no torque to the Earth if the Earth were perfectly spherical, but it is not. The Earth bulges slightly at the Equator. As indicated in Figure 26, the effect of the Sun's gravity on the near bulge (larger than it is on the far bulge) results in a net torque about the centre of the Earth. When the Earth is on the other side of the Sun, the net torque remains in the same direction. The torque

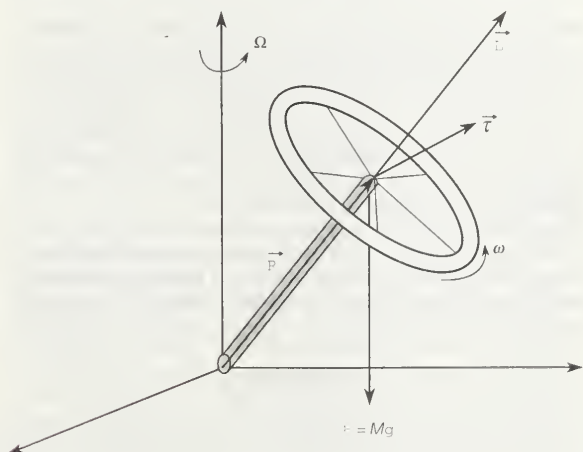


Figure 25: A gyroscope (see text).

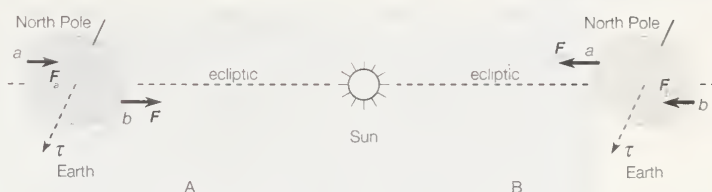


Figure 26: Forces acting on equatorial bulges in (A) the summer and (B) the winter cause the axis of the Earth to precess (see text).

After R.P. Olenick, T.M. Apostol, and D.L. Goodstein, *The Mechanical Universe: Introduction to Mechanics and Heat* (1985), Cambridge University Press

is small but persistent. It causes the axis of the Earth to precess, about one revolution every 25,800 years.

As seen from the Earth, the Sun passes through the plane of the Equator twice each year. These points are called the equinoxes, and on the days of the equinoxes the hours of daylight and night are equal. From antiquity it has been known that the point in the sky where the Sun intersects the plane of the Equator is not the same each year but rather drifts very slowly to the west. This ancient observation, first explained by Newton, is due to the precession of the Earth's axis.

Precession of the equinoxes

Analytic approaches

Classical mechanics can, in essence, be reduced to Newton's laws, starting with the second law, in the form

$$F = \frac{dp}{dt}. \tag{91}$$

If the net force acting on a particle is F , knowledge of F permits the momentum p to be found; and knowledge of p permits the position r to be found, by solving the equation

$$\frac{dr}{dt} = \frac{p}{m}. \tag{92}$$

These solutions give the components of p —that is, p_x , p_y , and p_z —and the components of r — x , y , and z —each as a function of time. To complete the solution, the value of each quantity— p_x , p_y , p_z , x , y , and z —must be known at some definite time, say, $t = 0$. If there is more than one particle, an equation in the form of equation (91) must be written for each particle, and the solution will involve finding the six variables x , y , z , p_x , p_y , and p_z , for each particle as a function of time, each once again subject to some initial condition. The equations may not be independent, however. For example, if the particles interact with one another, the forces will be related by Newton's third law. In this case (and others), the forces may also depend on time.

If the problem involves more than a very few particles, this method of solution quickly becomes intractable. Furthermore, in many cases it is not useful to express the problem purely in terms of particles and forces. Consider, for example, the problem of a sphere or cylinder rolling without slipping on a plane surface. Rolling without slipping is produced by friction due to forces acting between atoms in the rolling body and atoms in the plane, but the interactions are very complex; they probably are not fully understood even today, and one would like to be able to formulate and solve the problem without introducing them or needing to understand them. For all these reasons, methods that go beyond solving equations (91) and (92) have had to be introduced into classical mechanics.

The methods that have been introduced do not involve new physics. In fact, they are deduced directly from Newton's laws. They do, however, involve new concepts, new language to describe those concepts, and the adoption of powerful mathematical techniques. Some of those methods are briefly surveyed here.

CONFIGURATION SPACE

The position of a single particle is specified by giving its three coordinates, x , y , and z . To specify the positions of two particles, six coordinates are needed, x_1 , y_1 , z_1 , x_2 , y_2 , z_2 . If there are N particles, $3N$ coordinates will be needed. Imagine a system of $3N$ mutually orthogonal co-

ordinates in a $3N$ -dimensional space (a space of more than three dimensions is a purely mathematical construction, sometimes known as a hyperspace). To specify the exact position of one single point in this space, $3N$ coordinates are needed. However, one single point can represent the entire configuration of all N particles in the problem. Furthermore, the path of that single point as a function of time is the complete solution of the problem. This $3N$ -dimensional space is called configuration space.

Problem constraints

Configuration space is particularly useful for describing what is known as constraints on a problem. Constraints are generally ways of describing the effects of forces that are best not explicitly introduced into the problem. For example, consider the simple case of a falling body near the surface of the Earth. The equations of motion—equations (4), (5), and (6)—are valid only until the body hits the ground. Physically, this restriction is due to forces between atoms in the falling body and atoms in the ground, but, as a practical matter, it is preferable to say that the solutions are valid only for $z > 0$ (where $z = 0$ is ground level). This constraint, in the form of an inequality, is very difficult to incorporate directly into the equations of the problem. In the language of configuration space, however, one merely needs to specify that the problem is being solved only in the region of configuration space for which $z > 0$.

Notice that the constraint mentioned above, rolling without sliding on a plane, cannot easily be described in configuration space, since it is basically a condition on relative velocities of rotation and translation; but another constraint, that the body is restricted to motion along the plane, is easily described in configuration space.

Another type of constraint specifies that a body is rigid. Then, even though the body is composed of a very large number of atoms, it is not necessary to find separately the x , y , and z coordinate of each atom because these are related to those of the other atoms by the condition of rigidity. A careful analysis yields that, rather than needing $3N$ coordinates (where N may be, for example, 10^{24} atoms) only 6 are needed: 3 to specify the position of the centre of mass and 3 to give the orientation of the body. Thus, in this case, the constraint has reduced the number of independent coordinates from $3N$ to 6. Rather than restricting the behaviour of the system to a portion of the original $3N$ -dimensional configuration space, it is possible to describe the system in a much simpler 6-dimensional configuration space. It should be noted, however, that the six coordinates are not necessarily all distances. In fact, the most convenient coordinates are three distances (the x , y , and z coordinates of the centre of mass of the body) and three angles, which specify the orientation of a set of axes fixed in the body relative to a set of axes fixed in space. This is an example of the use of constraints to reduce the number of dynamic variables in a problem (the x , y , and z coordinates of each particle) to a smaller number of generalized dynamic variables, which need not even have the same dimensions as the original ones.

THE PRINCIPLE OF VIRTUAL WORK

A special class of problems in mechanics involves systems in equilibrium. The problem is to find the configuration of the system, subject to whatever constraints there may be, when all forces are balanced. The body or system will be at rest (in the inertial rest frame of its centre of mass), meaning that it occupies one point in configuration space for all time. The problem is to find that point. One criterion for finding that point, which makes use of the calculus of variations, is called the principle of virtual work.

Principle of virtual work

According to the principle of virtual work, any infinitesimal virtual displacement in configuration space, consistent with the constraints, requires no work. A virtual displacement means an instantaneous change in coordinates (a real displacement would require finite time during which particles might move and forces might change). To express the principle, label the generalized coordinates $r_1, r_2, \dots, r_n, \dots$. Then if F_i is the net component of generalized force acting along the coordinate r_i ,

$$\sum_i F_i dr_i = 0. \tag{93}$$

Here, $F_i dr_i$ is the work done when the generalized coordinate is changed by the infinitesimal amount dr_i . If r_i is a real coordinate (say, the x coordinate of a particle), then F_i is a real force. If r_i is a generalized coordinate (say, an angular displacement of a rigid body), then F_i is the generalized force such that $F_i dr_i$ is the work done (for an angular displacement, F_i is a component of torque).

Take two simple examples to illustrate the principle. First consider two particles that are restricted to motion in the x direction and are constrained by a taut string connecting them. If their x coordinates are called x_1 and x_2 , then $F_1 dx_1 + F_2 dx_2 = 0$ according to the principle of virtual work. But the taut string requires that the particles be displaced the same amount, so that $dx_1 = dx_2$, with the result that $F_1 + F_2 = 0$. The particles might be in equilibrium, for example, under equal and opposite forces, but F_1 and F_2 do not need individually to be zero. This is generally true of the F_i in equation (93). As a second example, consider a rigid body in space. Here, the constraint of rigidity has already been expressed by reducing the coordinate space to that of six generalized coordinates. These six coordinates (x , y , z , and three angles) can change quite independently of one another. In other words, in equation (93), the six dr_i are arbitrary. Thus, the only way equation (93) can be satisfied is if all six F_i are zero. This means that the rigid body can have no net component of force and no net component of torque acting on it. Of course, this same conclusion was reached earlier (see *Statics*) by less abstract arguments.

LAGRANGE'S AND HAMILTON'S EQUATIONS

Elegant and powerful methods have also been devised for solving dynamic problems with constraints. One of the best known is called Lagrange's equations. The Lagrangian L is defined as $L = T - V$, where T is the kinetic energy and V the potential energy of the system in question. Generally speaking, the potential energy of a system depends on the coordinates of all its particles; this may be written as $V = V(x_1, y_1, z_1, x_2, y_2, z_2, \dots)$. The kinetic energy generally depends on the velocities, which, using the notation $v_x = dx/dt = \dot{x}$, may be written $T = T(\dot{x}_1, \dot{y}_1, \dot{z}_1, \dot{x}_2, \dot{y}_2, \dot{z}_2, \dots)$. Thus, a dynamic problem has six dynamic variables for each particle—that is, x , y , z and \dot{x} , \dot{y} , \dot{z} —and the Lagrangian depends on all $6N$ variables if there are N particles.

In many problems, however, the constraints of the problem permit equations to be written relating at least some of these variables. In these cases, the $6N$ related dynamic variables may be reduced to a smaller number of independent generalized coordinates (written symbolically as $q_1, q_2, \dots, q_n, \dots$) and generalized velocities (written as $\dot{q}_1, \dot{q}_2, \dots, \dot{q}_n, \dots$), just as, for the rigid body, $3N$ coordinates were reduced to six independent coordinates (each of which has an associated velocity). The Lagrangian, then, may be expressed as a function of all the q_i and \dot{q}_i . It is possible, starting from Newton's laws only, to derive Lagrange's equations

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} - \frac{\partial L}{\partial q_i} = 0, \tag{94}$$

where the notation $\partial L / \partial q_i$ means differentiate L with respect to q_i only, holding all other variables constant. There is one equation of the form (94) for each of the generalized coordinates q_i (e.g., six equations for a rigid body), and their solutions yield the complete dynamics of the system. The use of generalized coordinates allows many coupled equations of the form (91) to be reduced to fewer, independent equations of the form (94).

There is an even more powerful method called Hamilton's equations. It begins by defining a generalized momentum p_i , which is related to the Lagrangian and the generalized velocity \dot{q}_i by $p_i = \partial L / \partial \dot{q}_i$. A new function, the Hamiltonian, is then defined by $H = \sum_i \dot{q}_i p_i - L$. From this point it is not difficult to derive

$$\dot{q}_i = \frac{\partial H}{\partial p_i} \tag{95}$$

and

Use of generalized coordinates and velocities

$$-\dot{p}_i = \frac{\partial H}{\partial q_i} \quad (96)$$

These are called Hamilton's equations. There are two of them for each generalized coordinate. They may be used in place of Lagrange's equations, with the advantage that only first derivatives—not second derivatives—are involved.

The Hamiltonian method is particularly important because of its utility in formulating quantum mechanics. However, it is also significant in classical mechanics. If

the constraints in the problem do not depend explicitly on time, then it may be shown that $H = T + V$, where T is the kinetic energy and V is the potential energy of the system—*i.e.*, the Hamiltonian is equal to the total energy of the system. Furthermore, if the problem is isotropic (H does not depend on direction in space) and homogeneous (H does not change with uniform translation in space), then Hamilton's equations immediately yield the laws of conservation of angular momentum and linear momentum, respectively. (D.L.G.)

CELESTIAL MECHANICS

Celestial mechanics, in the broadest sense, is the application of classical mechanics to the motion of celestial bodies acted on by any of several types of forces. By far the most important force experienced by these bodies, and much of the time the only important force, is that of their mutual gravitational attraction. But other forces can be important as well, such as atmospheric drag on artificial satellites, the pressure of radiation on dust particles, and even electromagnetic forces on dust particles if they are electrically charged and moving in a magnetic field. The term celestial mechanics is sometimes assumed to refer only to the analysis developed for the motion of point mass particles moving under their mutual gravitational attractions, with emphasis on the general orbital motions of solar system bodies. The term astrodynamics is often used to refer to the celestial mechanics of artificial satellite motion. Dynamic astronomy is a much broader term, which, in addition to celestial mechanics and astrodynamics, is usually interpreted to include all aspects of celestial body motion (*e.g.*, rotation, tidal evolution, mass and mass distribution determinations for stars and galaxies, fluid motions in nebulae, and so forth).

Dynamic
astronomy

HISTORICAL BACKGROUND

Early theories. Celestial mechanics has its beginnings in early astronomy in which the motions of the Sun, the Moon, and the five planets visible to the unaided eye—Mercury, Venus, Mars, Jupiter, and Saturn—were observed and analyzed. The word planet is derived from the Greek word for wanderer, and it was natural for some cultures to elevate these objects moving against the fixed background of the sky to the status of gods; this status survives in some sense today in astrology, where the positions of the planets and Sun are thought to somehow influence the lives of individuals on Earth. The divine status of the planets and their supposed influence on human activities may have been the primary motivation for careful, continued observations of planetary motions and for the development of elaborate schemes for predicting their positions in the future.

The Greek astronomer Ptolemy (who lived in Alexandria about AD 140) proposed a system of planetary motion in which the Earth was fixed at the centre and all the planets, the Moon, and the Sun orbited around it. As seen by an observer on the Earth, the planets move across the sky at a variable rate. They even reverse their direction of motion occasionally but resume the dominant direction of motion after a while. To describe this variable motion, Ptolemy assumed that the planets revolved around small circles called epicycles at a uniform rate while the centre of the epicyclic circle orbited the Earth on a large circle called a deferent. Other variations in the motion were accounted for by offsetting the centres of the deferent for each planet from the Earth by a short distance. By choosing the combination of speeds and distances appropriately, Ptolemy was able to predict the motions of the planets with considerable accuracy. His scheme was adopted as absolute dogma and survived more than 1,000 years until the time of Copernicus.

Copernicus'
heliocentric
model

Nicolaus Copernicus assumed that the Earth was just another planet that orbited the Sun along with the other planets. He showed that this heliocentric (centred on the Sun) model was consistent with all observations and that it was far simpler than Ptolemy's scheme. His belief that

planetary motion had to be a combination of uniform circular motions forced him to include a series of epicycles to match the motions in the noncircular orbits. The epicycles were like terms in the Fourier series that are used to represent planetary motions today. (A Fourier series is an infinite sum of periodic terms that oscillate between positive and negative values in a smooth way, where the frequency of oscillation changes from term to term. They represent better and better approximations to other functions as more and more terms are kept.) Copernicus also determined the relative scale of his heliocentric solar system, with results that are remarkably close to the modern determination.

Tycho Brahe (1546–1601), who was born three years after Copernicus' death and three years after the publication of the latter's heliocentric model of the solar system, still embraced a geocentric model, but he had only the Sun and the Moon orbiting the Earth and all the other planets orbiting the Sun. Although this model is mathematically equivalent to the heliocentric model of Copernicus, it represents an unnecessary complication and is physically incorrect. Tycho's greatest contribution was the more than 20 years of celestial observations he collected; his measurements of the positions of the planets and stars had an unprecedented accuracy of approximately 2 arc minutes. (An arc minute is $1/60$ of a degree.)

Kepler's laws of planetary motion. Tycho's observations were inherited by Johannes Kepler (1571–1630), who was employed by Tycho shortly before the latter's death. From these precise positions of the planets at correspondingly accurate times, Kepler empirically determined his famous three laws describing planetary motion: (1) the orbits of the planets are ellipses with the Sun at one focus; (2) the radial line from the Sun to the planet sweeps out equal areas in equal times; and (3) the ratio of the squares of the periods of revolution around the Sun of any two planets equal the ratio of the cubes of the semimajor axes of their respective orbital ellipses.

An ellipse (Figure 27) is a plane curve defined such that the sum of the distances from any point G on the ellipse to two fixed points (S and S' in Figure 27) is constant. The two points S and S' are called foci, and the straight line on which these points lie between the extremes of the ellipse at A and P is referred to as the major axis of the ellipse. Hence, $GS + GS' = AP = 2a$ in Figure 27, where a is the semimajor axis of the ellipse. A focus is separated from the centre C of the ellipse by the fractional part of the semimajor axis given by the product ae , where $e < 1$ is called the eccentricity. Thus, $e = 0$ corresponds to a circle. If the Sun is at the focus S of the ellipse, the point P at which the planet is closest to the Sun is called the perihelion, and the most distant point in the orbit A is the aphelion. The term helion refers specifically to the Sun as the primary body about which the planet is orbiting. As the points P and A are also called apsides, periapsis and apoapsis are often used to designate the corresponding points in an orbit about any primary body, although more specific terms, such as perigee and apogee for the Earth, are often used to indicate the primary body. If G is the instantaneous location of a planet in its orbit, the angle f , called the true anomaly, locates this point relative to the perihelion P with the Sun (or focus S) as the origin, or vertex, of the angle. The angle u , called the eccentric anomaly, also locates G relative to P but with the centre

of the ellipse as the origin rather than the focus S . An angle called the mean anomaly l (not shown in Figure 27) is also measured from P with S as the origin; it is defined to increase uniformly with time and to equal the true anomaly f at perihelion and aphelion.

Law of
equal areas

Kepler's second law is also illustrated in Figure 27. If the time required for the planet to move from P to F is the same as that to move from D to E , the areas of the two shaded regions will be equal according to the second law. The validity of the second law means a planet must have a higher than average velocity near perihelion and a lower than average velocity near aphelion. The angular velocity (rate of change of the angle f) must vary around the orbit in a similar way. The average angular velocity, called the mean motion, is the rate of change of the mean anomaly l defined above.

The third law can be used to determine the distance of a planet from the Sun if one knows its orbital period, or vice versa. In particular, if time is measured in years and distance in units of the semimajor axis of the Earth's orbit (*i.e.*, the mean distance of the Earth to the Sun, known as an astronomical unit, or AU), the third law can be written $\tau^2 = a^3$, where τ is the orbital period.

Newton's laws of motion. The empirical laws of Kepler describe planetary motion, but Kepler made no attempt to define or constrain the underlying physical processes governing the motion. It was Isaac Newton who accomplished that feat in the late 17th century. Newton defined momentum as being proportional to velocity with the constant of proportionality being defined as mass. (As described earlier, momentum is a vector quantity in the sense that the direction of motion as well as the magnitude is included in the definition.) Newton then defined force (also a vector quantity) in terms of its effect on moving objects and in the process formulated his three laws of motion: (1) The momentum of an object is constant unless an outside force acts on the object; this means that any object either remains at rest or continues uniform motion in a straight line unless acted on by a force. (2) The time rate of change of the momentum of an object is equal to the force acting on the object. (3) For every action (force) there is an equal and opposite reaction (force). The first law is seen to be a special case of the second law. Galileo, the great Italian contemporary of Kepler who adopted the Copernican point of view and promoted it vigorously, anticipated Newton's first two laws with his experiments in mechanics. But it was Newton who defined them precisely, established the basis of classical mechanics, and set the stage for its application as celestial mechanics to the motions of bodies in space.

According to the second law, a force must be acting on a planet to cause its path to curve toward the Sun. Newton and others noted that the acceleration of a body in uni-

form circular motion must be directed toward the centre of the circle; furthermore, if several objects were in circular motion around the same centre at various separations r and their periods of revolution varied as $r^{3/2}$, as Kepler's third law indicated for the planets, then the acceleration—and thus, by Newton's second law, the force as well—must vary as $1/r^2$. By assuming this attractive force between point masses, Newton showed that a spherically symmetric mass distribution attracted a second body outside the sphere as if all the spherically distributed mass were contained in a point at the centre of the sphere. Thus the attraction of the planets by the Sun was the same as the gravitational force attracting objects to the Earth. Newton further concluded that the force of attraction between two massive bodies was proportional to the inverse square of their separation and to the product of their masses, known as the law of universal gravitation. Kepler's laws are derivable from Newton's laws of motion with a central force of gravity varying as $1/r^2$ from a fixed point, and Newton's law of gravity is derivable from Kepler's laws if one assumes Newton's laws of motion.

In addition to formulating the laws of motion and of gravity, Newton also showed that a point mass moving about a fixed centre of force, which varies as the inverse square of the distance away from the centre, follows an elliptical path if the initial velocity is not too large, a hyperbolic path for high initial velocities, and a parabolic path for intermediate velocities. In other words, a sequence of orbits in Figure 27 with the perihelion distance SP fixed but with the velocity at P increasing from orbit to orbit is characterized by a corresponding increase in the orbital eccentricity e from orbit to orbit such that $e < 1$ for bound elliptical orbits, $e = 1$ for a parabolic orbit, and $e > 1$ for a hyperbolic orbit. Many comets have nearly parabolic orbits for their first pass into the inner solar system, whereas spacecraft may have nearly hyperbolic orbits relative to a planet they are flying by while they are close to the planet.

Orbital
eccentricity

Throughout history, the motion of the planets in the solar system has served as a laboratory to constrain and guide the development of celestial mechanics in particular and classical mechanics in general. In modern times, increasingly precise observations of celestial bodies have been matched by increasingly precise predictions for future positions—a combination that became a test for Newton's law of gravitation itself. Although the lunar motion (within observational errors) seemed consistent with a gravitational attraction between point masses that decreased exactly as $1/r^2$, this law of gravitation was ultimately shown to be an approximation of the more complete description of gravity given by the theory of general relativity. Similarly, a discrepancy of roughly 40 arc seconds per century between the observed rate of advance of Mercury's perihelion and that predicted by planetary perturbations with Newtonian gravity is almost precisely accounted for with Einstein's general theory of relativity. That this small discrepancy could be confidently asserted as real was a triumph of quantitative celestial mechanics.

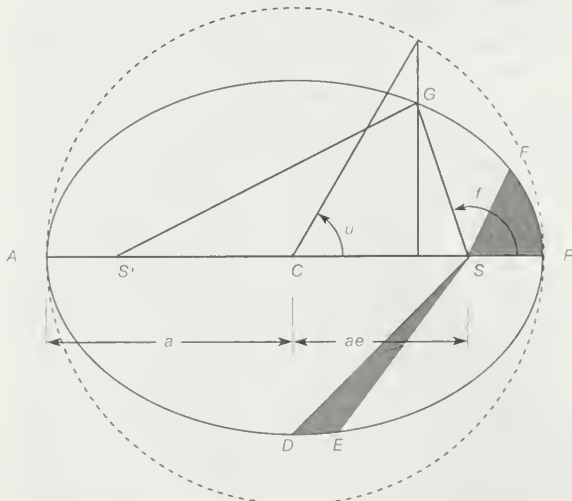


Figure 27: The orbital elements a (the semimajor axis) and e (the eccentricity) characterize an elliptical orbit; the angles f and u allow location of the position of a planet on the orbit relative to the point P ; the shaded areas illustrate Kepler's second law (see text).

PERTURBATIONS AND PROBLEMS OF TWO BODIES

The approximate nature of Kepler's laws. The constraints placed on the force for Kepler's laws to be derivable from Newton's laws were that the force must be directed toward a central fixed point and that the force must decrease as the inverse square of the distance. In actuality, however, the Sun, which serves as the source of the major force, is not fixed but experiences small accelerations because of the planets, in accordance with Newton's second and third laws. Furthermore, the planets attract one another, so that the total force on a planet is not just that due to the Sun; other planets perturb the elliptical motion that would have occurred for a particular planet if that planet had been the only one orbiting an isolated Sun. Kepler's laws therefore are only approximate. The motion of the Sun itself means that, even when the attractions by other planets are neglected, Kepler's third law must be replaced by $(M + m_i)\tau^2 \propto a^3$, where m_i is one of the planetary masses and M is the Sun's mass. That Kepler's laws are such good approximations to the actual planetary motions results from the fact that all the planetary masses are very

small compared to that of the Sun. The perturbations of the elliptic motion are therefore small, and the coefficient $M + m_i \approx M$ for all the planetary masses m_i , means that Kepler's third law is very close to being true.

Newton's second law for a particular mass is a second-order differential equation that must be solved for whatever forces may act on the body if its position as a function of time is to be deduced. The exact solution of this equation, which resulted in a derived trajectory that was an ellipse, parabola, or hyperbola, depended on the assumption that there were only two point particles interacting by the inverse square force. Hence, this "gravitational two-body problem" has an exact solution that reproduces Kepler's laws. If one or more additional bodies also interact with the original pair through their mutual gravitational interactions, no exact solution for the differential equations of motion of any of the bodies involved can be obtained. As was noted above, however, the motion of a planet is almost elliptical, since all masses involved are small compared to the Sun. It is then convenient to treat the motion of a particular planet as slightly perturbed elliptical motion and to determine the changes in the parameters of the ellipse that result from the small forces as time progresses. It is the elaborate developments of various perturbation theories and their applications to approximate the exact motions of celestial bodies that has occupied celestial mechanics since Newton's time.

Perturbations of elliptical motion. So far the following orbital parameters, or elements, have been used to describe elliptical motion: the orbital semimajor axis a ; the orbital eccentricity e , and, to specify position in the orbit relative to the perihelion, either the true anomaly f , the eccentric anomaly u , or the mean anomaly l . Three more orbital elements are necessary to orient the ellipse in space, since that orientation will change because of the perturbations. The most commonly chosen of these additional parameters are illustrated in Figure 28, where the reference plane is chosen arbitrarily to be the plane of the ecliptic, which is the plane of the Earth's orbit defined by the path of the Sun on the sky. (For motion of a near-Earth artificial satellite, the most convenient reference plane would be that of the Earth's Equator.) Angle i is the inclination of the orbital plane to the reference plane. The line of nodes is the intersection of the orbit plane with the reference plane, and the ascending node is that point where the planet travels from below the reference plane (south) to above the reference plane (north). The ascending node is described by its angular position measured from a reference point on the ecliptic plane, such as the vernal equinox; the angle Ω is called the longitude of the ascending node. Angle ω (called the argument of perihelion) is the angular distance from the ascending node to the perihelion measured in the orbit plane.

Orbital
elements

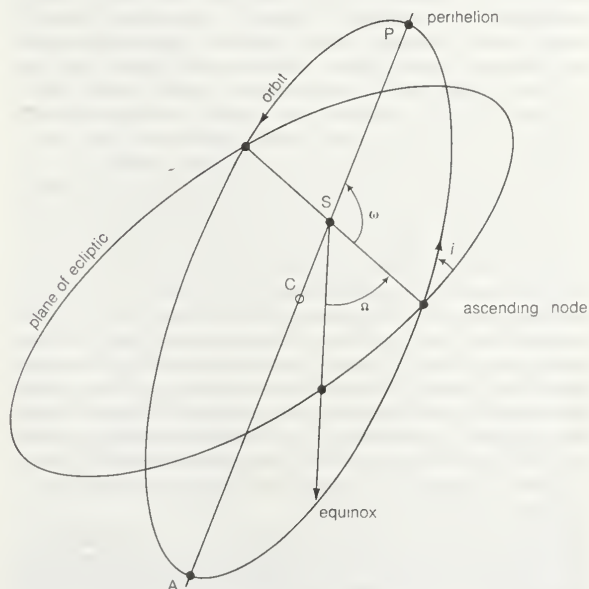


Figure 28: Orbital elements Ω , ω , and i orienting the ellipse.

For the two-body problem, all the orbital parameters a , e , i , Ω , and ω are constants. A sixth constant T , the time of perihelion passage (*i.e.*, any date at which the object in orbit was known to be at perihelion), may be used to replace f , u , or l , and the position of the planet in its fixed elliptic orbit can be determined uniquely at subsequent times. These six constants are determined uniquely by the six initial conditions of three components of the position vector and three components of the velocity vector relative to a coordinate system that is fixed with respect to the reference plane. When small perturbations are taken into account, it is convenient to consider the orbit as an instantaneous ellipse whose parameters are defined by the instantaneous values of the position and velocity vectors, since for small perturbations the orbit is approximately an ellipse. In fact, however, perturbations cause the six formerly constant parameters to vary slowly, and the instantaneous perturbed orbit is called an osculating ellipse; that is, the osculating ellipse is that elliptical orbit that would be assumed by the body if all the perturbing forces were suddenly turned off.

First-order differential equations describing the variation of the six orbital parameters can be constructed for each planet or other celestial body from the second-order differential equations that result by equating the mass times the acceleration of a body to the sum of all the forces acting on the body (Newton's second law). These equations are sometimes called the Lagrange planetary equations after their derivation by the great Italian-French mathematician Joseph-Louis Lagrange (1736–1813). As long as the forces are conservative and do not depend on the velocities—*i.e.*, there is no loss of mechanical energy through such processes as friction—they can be derived from partial derivatives of a function of the spatial coordinates only, called the potential energy, whose magnitude depends on the relative separations of the masses.

In the case where all the forces are derivable from such potential energy, the total energy of a system of any number of particles—*i.e.*, the kinetic energy plus the potential energy—is constant. The kinetic energy of a single particle is one-half its mass times the square of its velocity, and the total kinetic energy is the sum of such expressions for all the particles being considered. The conservation of energy principle is thus expressed by an equation relating the velocities of all the masses to their positions at any time. The partial derivatives of the potential energy with respect to spatial coordinates are transformed into particle derivatives of a disturbing function with respect to the orbital elements in the Lagrange equations, where the disturbing function vanishes if all bodies perturbing the elliptic motion are removed. Like Newton's equations of motion, Lagrange's differential equations are exact, but they can be solved only numerically on a computer or analytically by successive approximations. In the latter process, the disturbing function is represented by a Fourier series, with convergence of the series (successive decrease in size and importance of the terms) depending on the size of the orbital eccentricities and inclinations. Clever changes of variables and other mathematical tricks are used to increase the time span over which the solutions (also represented by series) are good approximations to the real motion. These series solutions usually diverge, but they still represent the actual motions remarkably well for limited periods of time. One of the major triumphs of celestial mechanics using these perturbation techniques was the discovery of Neptune in 1846 from its perturbations of the motion of Uranus.

Disturbing
function

Examples of perturbations. Some of the variations in the orbital parameters caused by perturbations can be understood in simple terms. The lunar orbit is inclined to the ecliptic plane by about 5° , and the longitude of its ascending node on the ecliptic plane (Ω in Figure 28) is observed to regress (Ω decreasing) a complete revolution in 18.61 years. The Sun is the dominant cause of this regression of the lunar node. When the Moon is closer to the Sun than the Earth, the Sun accelerates the Moon slightly more than it accelerates the Earth. This difference in the accelerations is what perturbs the lunar motion around the Earth. The Moon does not fly off in this situation since

the acceleration of the Moon toward the Earth is much larger than the difference between the Sun's accelerations of the Earth and the Moon.

The Sun, of course, is always in the ecliptic plane, since its apparent path among the stars defines the plane. This means that the perturbing acceleration just defined will always be pointed slightly toward the ecliptic plane whenever the Moon is below or above this plane in its orbital motion about the Earth. This tendency to pull the Moon toward the ecliptic plane means that the Moon will cross the plane on each half orbit at a longitude that is slightly smaller than the longitude at which it would have crossed if the Sun had not been there. Thus, the line of nodes will have regressed. The instantaneous rate at which the node regresses varies as the geometry changes during the Moon's motion around the Earth, and during the Earth-Moon system's motion around the Sun, but there is always a net regression. Such a change that is always in the same direction as time increases is called a secular perturbation. Superposed on the secular perturbation of the longitude of the node are periodic perturbations (periodically changing their direction), which are revealed by the fact that the rate of secular regression of the node is not constant in time. The Sun causes a secular increase in the longitude of the lunar perigee ($\Omega + \omega$ in Figure 28) of one complete revolution in 8.85 years, as well as periodic perturbations in the inclination, eccentricity, and mean motion.

For near-Earth artificial satellites, the deviation of the Earth's mass distribution from spherical symmetry is the dominant cause of the perturbations from pure elliptic motion. The most important deviation is the equatorial bulge of the Earth due to its rotation. If, for example, the Earth were a sphere with a ring of mass around its Equator, the ring would give to a satellite whose orbit is inclined to the Equator a component of acceleration toward the Equator plane whenever the satellite was above or below this plane. By an argument similar to that for the Moon acted on by the Sun, this acceleration would cause the line of nodes of a close satellite orbit to regress a little more than 5° per day.

As a final example, the distribution of continents and oceans and the varying mass densities in the Earth's mantle (the layer underlying the crust) lead to a slight deviation of the Earth's gravitational force field from axial symmetry. Usually this causes only short-period perturbations of low amplitude for near-Earth satellites. However, communications or weather satellites that are meant to maintain a fixed longitude over the Equator (*i.e.*, geostationary satellites, which orbit synchronously with the Earth's rotation) are destabilized by this deviation except at two longitudes. If the axial asymmetry is represented by a slightly elliptical Equator, the difference between the major and minor axis of the ellipse is about 64 metres, with the major axis located about 35° W. A satellite at a position slightly ahead of the long axis of the elliptical Equator will experience a component of acceleration opposite its direction of orbital motion (as if a large mountain were pulling it back). This acceleration makes the satellite fall closer to the Earth and increases its mean motion, causing it to drift further ahead of the axial bulge on the Equator. If the satellite is slightly behind the axial bulge, it experiences an acceleration in the direction of its motion. This makes the satellite move away from the Earth with a decrease in its mean motion, so that it will drift further behind the axial bulge. The synchronous Earth satellites are thus repelled from the long axis of the equatorial ellipse and attracted to the short axis, and compensating accelerations, usually from onboard jets, are required to stabilize a satellite at any longitude other than the two corresponding to the ends of the short axis of the axial bulge. (The jets are actually required for any longitude, as they must also compensate for other perturbations such as radiation pressure.)

THE THREE-BODY PROBLEM

The inclusion of solar perturbations of the motion of the Moon results in a "three-body problem" (Earth-Moon-Sun), which is the simplest complication of the completely solvable two-body problem discussed above. When the Earth, Moon, and Sun are considered to be point masses,

this particular three-body problem is called "the main problem of the lunar theory," which has been studied extensively with a variety of methods beginning with Newton. Although the three-body problem has no complete analytic solution in closed form, various series solutions by successive approximations achieve such accuracy that complete theories of the lunar motion must include the effects of the nonspherical mass distributions of both the Earth and the Moon as well as the effects of the planets if the precision of the predicted positions is to approach that of the observations. Most of the schemes for the main problem are partially numerical and therefore apply only to the lunar motion. An exception is the completely analytic work of the French astronomer Charles-Eugène Delaunay (1816–72), who exploited and developed the most elegant techniques of classical mechanics pioneered by his contemporary, the Irish astronomer and mathematician William R. Hamilton (1805–65). Delaunay could predict the position of the Moon to within its own diameter over a 20-year time span. Since his development was entirely analytic, the work was applicable to the motions of satellites about other planets where the series expansions converged much more rapidly than they did for the application to the lunar motion.

Delaunay's work on the lunar theory demonstrates some of the influence that celestial mechanics has had on the development of the techniques of classical mechanics. This close link between the development of classical mechanics and its application to celestial mechanics was probably no better demonstrated than in the work of the French mathematician Henri Poincaré (1854–1912). Poincaré, along with other great mathematicians such as George D. Birkhoff (1884–1944), Aurel Wintner (1903–58), and Andrey N. Kolmogorov (1903–87), placed celestial mechanics on a more sound mathematical basis and was less concerned with quantitatively accurate prediction of celestial body motion. Poincaré demonstrated that the series solutions in use in celestial mechanics for so long generally did not converge but that they could give accurate descriptions of the motion for significant periods of time in truncated form. The elaborate theoretical developments in celestial and classical mechanics have received more attention recently with the realization that a large class of motions are of an irregular or chaotic nature and require fundamentally different approaches for their description.

The restricted three-body problem. The simplest form of the three-body problem is called the restricted three-body problem, in which a particle of infinitesimal mass moves in the gravitational field of two massive bodies orbiting according to the exact solution of the two-body problem. (The particle with infinitesimal mass, sometimes called a massless particle, does not perturb the motions of the two massive bodies.) There is an enormous literature devoted to this problem, including both analytic and numerical developments. The analytic work was devoted mostly to the circular, planar restricted three-body problem, where all particles are confined to a plane and the two finite masses are in circular orbits around their centre of mass (a point on the line between the two masses that is closer to the more massive). Numerical developments allowed consideration of the more general problem.

In the circular problem, the two finite masses are fixed in a coordinate system rotating at the orbital angular velocity, with the origin (axis of rotation) at the centre of mass of the two bodies. Lagrange showed that in this rotating frame there were five stationary points at which the massless particle would remain fixed if placed there. There are three such points lying on the line connecting the two finite masses: one between the masses and one outside each of the masses. The other two stationary points, called the triangular points, are located equidistant from the two finite masses at a distance equal to the finite mass separation. The two masses and the triangular stationary points are thus located at the vertices of equilateral triangles in the plane of the circular orbit.

There is a constant of the motion in the rotating frame that leads to an equation relating the velocity of the massless particle in this frame to its position. For given values of this constant it is possible to construct curves in the

Influence
on classical
mechanics

Syn-
chronous
Earth
satellites

Zero-velocity curves

plane on which the velocity vanishes. If such a zero-velocity curve is closed, the particle cannot escape from the interior of the closed zero-velocity curve if placed there with the constant of the motion equal to the value used to construct the curve. These zero-velocity curves can be used to show that the three collinear stationary points are all unstable in the sense that, if the particle is placed at one of these points, the slightest perturbation will cause it to move far away. The triangular points are stable if the ratio of the finite masses is less than 0.04, and the particle would execute small oscillations around one of the triangular points if it were pushed slightly away. Since the mass ratio of Jupiter to the Sun is about 0.001, the stability criterion is satisfied, and Lagrange predicted the presence of the Trojan asteroids at the triangular points of the Sun-Jupiter system 134 years before they were observed. Of course, the stability of the triangular points must also depend on the perturbations by any other bodies. Such perturbations are sufficiently small not to destabilize the Trojan asteroids. Single Trojan-like bodies have also been found orbiting at leading and trailing triangular points in the orbit of Saturn's satellite Tethys, at the leading triangular point in the orbit of another Saturnian satellite, Dione, and at the trailing point in the orbit of Mars.

Orbital resonances. There are stable configurations in the restricted three-body problem that are not stationary in the rotating frame. If, for example, Jupiter and the Sun are the two massive bodies, these stable configurations occur when the mean motions of Jupiter and the small particle—here an asteroid—are near a ratio of small integers. The orbital mean motions are then said to be nearly commensurate, and an asteroid that is trapped near such a mean motion commensurability is said to be in an orbital resonance with Jupiter. For example, the Trojan asteroids librate (oscillate) around the 1:1 orbital resonance (*i.e.*, the orbital period of Jupiter is in a 1:1 ratio with the orbital period of the Trojan asteroids); the asteroid Thule librates around the 4:3 orbital resonance; and several asteroids in the Hilda group librate around the 3:2 orbital resonance. There are several such stable orbital resonances among the satellites of the major planets and one involving the planets Pluto and Neptune. The analysis based on the restricted three-body problem cannot be used for the satellite resonances, however, except for the 4:3 resonance between Saturn's satellites Titan and Hyperion, since the participants in the satellite resonances usually have comparable masses.

Although the asteroid Griqua librates around the 2:1 resonance with Jupiter, and Alinda librates around the 3:1 resonance, the orbital commensurabilities 2:1, 7:3, 5:2, and 3:1 are characterized by an absence of asteroids in an otherwise rather highly populated, uniform distribution spanning all of the commensurabilities. These are the Kirkwood gaps in the distribution of asteroids, and the recent understanding of their creation and maintenance has introduced into celestial mechanics an entirely new concept of irregular, or chaotic, orbits in a system whose equations of motion are entirely deterministic.

Chaotic orbits. The French astronomer Michel Hénon and the American astronomer Carl Heiles discovered that when a system exhibiting periodic motion, such as a pendulum, is perturbed by an external force that is also periodic, some initial conditions lead to motions where the state of the system becomes essentially unpredictable (within some range of system states) at some time in the future, whereas initial conditions within some other set produce quasiperiodic or predictable behaviour. The unpredictable behaviour is called chaotic, and initial conditions that produce it are said to lie in a chaotic zone. If the chaotic zone is bounded, in the sense that only limited ranges of initial values of the variables describing the motion lead to chaotic behaviour, the uncertainty in the state of the system in the future is limited by the extent of the chaotic zone; that is, values of the variables in the distant future are completely uncertain only within those ranges of values within the chaotic zone. This complete uncertainty within the zone means the system will eventually come arbitrarily close to any set of values of the variables within the zone if given sufficient

time. Chaotic orbits were first realized in the asteroid belt.

A periodic term in the expansion of the disturbing function for a typical asteroid orbit becomes more important in influencing the motion of the asteroid if the frequency with which it changes sign is very small and its coefficient is relatively large. For asteroids orbiting near a mean motion commensurability with Jupiter, there are generally several terms in the disturbing function with large coefficients and small frequencies that are close but not identical. These "resonant" terms often dominate the perturbations of the asteroid motion so much that all the higher-frequency terms can be neglected in determining a first approximation to the perturbed motion. This neglect is equivalent to averaging the higher-frequency terms to zero; the low-frequency terms change only slightly during the averaging. If one of the frequencies vanishes on the average, the periodic term becomes nearly constant, or secular, and the asteroid is locked into an exact orbital resonance near the particular mean motion commensurability. The mean motions are not exactly commensurate in such a resonance, however, since the motion of the asteroid orbital node or perihelion is always involved (except for the 1:1 Trojan resonances).

For example, for the 3:1 commensurability, the angle $\theta = \lambda_A - 3\lambda_J + \varpi_A$ is the argument of one of the important periodic terms whose variation can vanish (zero frequency). Here $\lambda = \Omega + \omega + l$ is the mean longitude, the subscripts *A* and *J* refer to the asteroid and Jupiter, respectively, and $\varpi = \Omega + \omega$ is the longitude of perihelion (see Figure 28). Within resonance, the angle θ librates, or oscillates, around a constant value as would a pendulum around its equilibrium position at the bottom of its swing. The larger the amplitude of the equivalent pendulum, the larger its velocity at the bottom of its swing. If the velocity of the pendulum at the bottom of its swing, or, equivalently, the maximum rate of change of the angle θ , is sufficiently high, the pendulum will swing over the top of its support and be in a state of rotation instead of libration. The maximum value of the rate of change of θ for which θ remains an angle of libration (periodically reversing its variation) instead of one of rotation (increasing or decreasing monotonically) is defined as the half-width of the resonance.

Another term with nearly zero frequency when the asteroid is near the 3:1 commensurability has the argument $\theta' = \lambda_A - \lambda_J + 2\varpi_J$. The substitution of the longitude of Jupiter's perihelion for that of the asteroid means that the rates of change of θ and θ' will be slightly different. As the resonances are not separated much in frequency, there may exist values of the mean motion of the asteroid where both θ and θ' would be angles of libration if either resonance existed in the absence of the other. The resonances are said to overlap in this case, and the attempt by the system to librate simultaneously about both resonances for some initial conditions leads to chaotic orbital behaviour. The important characteristic of the chaotic zone for asteroid motion near a mean motion commensurability with Jupiter is that it includes a region where the asteroid's orbital eccentricity is large. During the variation of the elements over the entire chaotic zone as time increases, large eccentricities must occasionally be reached. For asteroids near the 3:1 commensurability with Jupiter, the orbit then crosses that of Mars, whose gravitational interaction in a close encounter can remove the asteroid from the 3:1 zone.

By numerically integrating many orbits whose initial conditions spanned the 3:1 Kirkwood gap region in the asteroid belt, Jack Wisdom, an American dynamicist who developed a powerful means of analyzing chaotic motions, found that the chaotic zone around this gap precisely matched the physical extent of the gap. There are no observable asteroids with orbits within the chaotic zone, but there are many just outside extremes of the zone. Preliminary work has indicated that the other Kirkwood gaps can be similarly accounted for. The realization that orbits governed by Newton's laws of motion and gravitation could have chaotic properties and that such properties could solve a long-standing problem in the celestial mechanics of the solar system is a major breakthrough in the subject.

Approximating asteroid motion

Kirkwood gaps

Chaotic zones

THE n -BODY PROBLEM

The general problem of n bodies, where n is greater than three, has been attacked vigorously with numerical techniques on powerful computers. Celestial mechanics in the solar system is ultimately an n -body problem, but the special configurations and relative smallness of the perturbations have allowed quite accurate descriptions of motions (valid for limited time periods) with various approximations and procedures without any attempt to solve the complete problem of n bodies. Examples are the restricted three-body problem to determine the effect of Jupiter's perturbations of the asteroids and the use of successive approximations of series solutions to sequentially add the effects of smaller and smaller perturbations for the motion of the Moon. In the general n -body problem, all bodies have arbitrary masses, initial velocities, and positions; the bodies interact through Newton's law of gravitation, and one attempts to determine the subsequent motion of all the bodies. Many numerical solutions for the motion of quite large numbers of gravitating particles have been successfully completed where the precise motion of individual particles is usually less important than the statistical behaviour of the group.

Numerical solutions. Numerical solutions of the exact equations of motion for n bodies can be formulated. Each body is subject to the gravitational attraction of all the others, and it may be subject to other forces as well. It is relatively easy to write the expression for the instantaneous acceleration (equation of motion) of each body if the position of all the other bodies is known, and expressions for all the other forces can be written (as they can for gravitational forces) in terms of the relative positions of the particles and other defining characteristics of the particle and its environment. Each particle is allowed to move under its instantaneous acceleration for a short time step. Its velocity and position are thereby changed, and the new values of the variables are used to calculate the acceleration for the next time step, and so forth. Of course, the real position and velocity of the particle after each time step will differ from the calculated values by errors of two types. One type results from the fact that the acceleration is not really constant over the time step, and the other from the rounding off or truncation of the numbers at every step of the calculation. The first type of error is decreased by taking shorter time steps. But this means more numerical operations must be carried out over a given span of time, and this increases the round-off error for a given precision of the numbers being carried in the calculation. The design of numerical algorithms, as well as the choice of precisions and step sizes that maximize the speed of the calculation while keeping the errors within reasonable bounds, is almost an art form developed by extensive experience and ingenuity. For example, a scheme exists for extrapolating the step size to zero in order to find the change in the variables over a relatively short time span, thereby minimizing the accumulation of error from this source. If the total energy of the system is theoretically conserved, its evaluation for values of the variables at the beginning and end of a calculation is a measure of the errors that have accumulated.

The motion of the nine planets of the solar system over time scales approaching its 4.6-billion-year age is a classic n -body problem, where $n = 10$ with the Sun included. The question of whether or not the solar system is ultimately stable—whether the current configuration of the planets will be maintained indefinitely under their mutual perturbations, or whether one or another planet will eventually be lost from the system or otherwise have its orbit drastically altered—is a long-standing one that might someday be answered through numerical calculation. The interplay of orbital resonances and chaotic orbits discussed above can be investigated numerically, and this interplay may be crucial in determining the stability of the solar system. Already it appears that the parameters defining the orbits of several planets, especially that of Pluto, vary over narrow chaotic zones, but whether or not this chaos can lead to instability if given enough time is still uncertain.

If accelerations are determined by summing all the pairwise interactions for the n particles, the computer time

per time step increases as n^2 . Practical computations for the direct calculation of the interactions between all the particles are thereby limited to $n < 10,000$. Therefore, for larger values of n , schemes are used where a particle is assumed to move in the force field of the remaining particles approximated by that due to a continuum mass distribution, or a "tree structure" is used where the effects of nearby particles are considered individually while larger and larger groups of particles are considered collectively as their distances increase. These later schemes have the capability of calculating the evolution of a very large system of particles using a reasonable amount of computer time with reasonable approximation. Values of n near 100,000 have been used in calculations determining the evolution of galaxies of stars. Also, the consequences for distribution of stars when two galaxies closely approach one another or even collide has been determined. Even calculations of the n -body problem where n changes with time have been completed in the study of the accumulation of larger bodies from smaller bodies via collisions in the process of the formation of the planets.

In all n -body calculations, very close approaches of two particles can result in accelerations so large and so rapidly changing that large errors are introduced or the calculation completely diverges. Accuracy can sometimes be maintained in such a close approach, but only at the expense of requiring very short time steps, which drastically slows the calculation. When n is small, as in some solar system calculations where two-body orbits still dominate, close approaches are sometimes handled by a change to a set of variables, usually involving the eccentric anomaly u , that vary much less rapidly during the encounter. In this process, called regularization, the encounter is traversed in less computer time while preserving reasonable accuracy. This process is impractical when n is large, so accelerations are usually artificially bounded on close approaches to prevent instabilities in the numerical calculation and to prevent slowing the calculation. For example, if several sets of particles were trapped in stable, close binary orbits, the very short time steps required to follow this rapid motion would bring the calculation to a virtual standstill, and such binary motion is not important in the overall evolution of, say, a galaxy of stars.

Algebraic maps. In numerical calculations for conservative systems with modest values of n over long time spans, such as those seeking a determination of the stability of the solar system, the direct solution of the differential equations governing the motions requires excessive time on any computer and accumulates excessive round-off error in the process. Excessive time also is required to explore thoroughly a complete range of orbital parameters in numerical experiments in order to determine the extent of chaotic zones in various configurations (e.g., those in the asteroid belt near orbital mean motion commensurabilities with Jupiter). A solution to this problem is the use of an algebraic map, which maps the space of system variables onto itself in such a way that the values of all the variables at one instant of time are expressed by algebraic relations in terms of the values of the variables at a fixed time in the past. The values at the next time step are determined by applying the same map to the values just obtained, and so on. The map is constructed by assuming that the motions of all the bodies are unperturbed for a given short time but are periodically "kicked" by the perturbing forces for only an instant. The continuous perturbations are thus replaced by periodic impulses. The values of the variables are "mapped" from one time step to the next by the fact that the unperturbed part of the motion is available from the exact solution of the two-body problem, and it is easy to solve the equations with all the perturbations over the short time of the impulse. Although this approximation does not produce exactly the same values of all the variables at some time in the future as those produced by a numerical solution of the differential equations starting with the same initial conditions, the qualitative behaviour is indistinguishable over long time periods. As computers can perform the algebraic calculations as much as 1,000 times faster than they can solve the corresponding differential equations, the computational time savings are

Regularization

Stability of the solar system

enormous and problems otherwise impossible to explore become tractable.

TIDAL EVOLUTION

This discussion has so far treated the celestial mechanics of bodies accelerated by conservative forces (total energy being conserved), including perturbations of elliptic motion by nonspherical mass distributions of finite-size bodies. However, the gravitational field of one body in close orbit about another will tidally distort the shape of the other body. Dissipation of part of the energy stored in these tidal distortions leads to a coupling that causes secular changes (always in the same direction) in the orbit and in the spins of both bodies. Since tidal dissipation accounts for the current spin states of several planets, the spin states of most of the planetary satellites and some of their orbital configurations, and the spins and orbits of close binary stars, it is appropriate that tides and their consequences be included in this discussion.

Tidal deformations of the Earth

The twice-daily high and low tides in the ocean are known by all who have lived near a coast. Few are aware, however, that the solid body of the Earth also experiences twice-daily tides with a maximum amplitude of about 30 centimetres. George Howard Darwin (1845–1912), the second son of Charles Darwin, the naturalist, was an astronomer-geophysicist who understood quantitatively the generation of the tides in the gravitational fields of tide-raising bodies, which are primarily the Moon and Sun for the Earth; he pointed out that the dissipation of tidal energy resulted in a slowing of the Earth's rotation while the Moon's orbit was gradually expanded. That any mass raises a tide on every other mass within its gravitational field follows from the fact that the gravitational force between two masses decreases as the inverse square of the distance between them.

In Figure 29, the accelerations due to mass m_s of three mass elements in the spherical mass m_p are proportional to the length of the arrows attached to each element. The element nearest m_s is accelerated more than the element at the centre of m_p and tends to leave the centre element behind; the element at the centre of m_p is accelerated more than the element farthest from m_s , and the latter tends to be left further behind. The point of view of a fictitious observer at the centre of m_p can be realized by subtracting the acceleration of the central mass element from that of each of the other two mass elements. If the mass elements were free, this observer would see the two extreme mass elements being accelerated in opposite directions away from his position at the centre, as illustrated in Figure 29B.

But the mass elements are not free; they are gravitationally attracted to one another and to the remaining mass in m_p . The gravitational acceleration of the mass elements on the surface of m_p toward the centre of m_p far exceeds

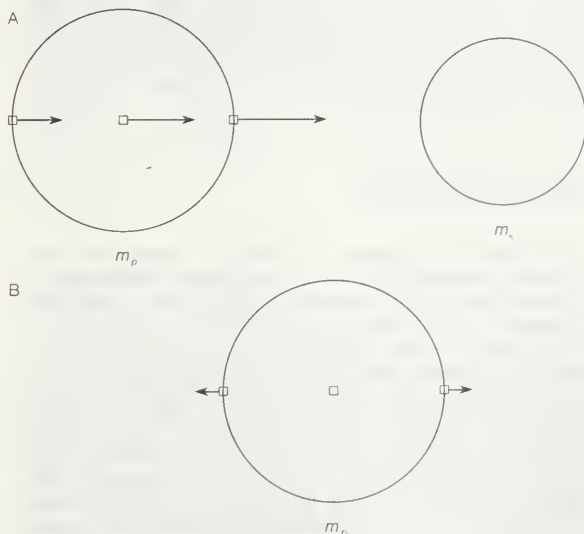


Figure 29: Variation of gravitational acceleration across a finite-sized body leading to differential acceleration relative to its centre (see text).

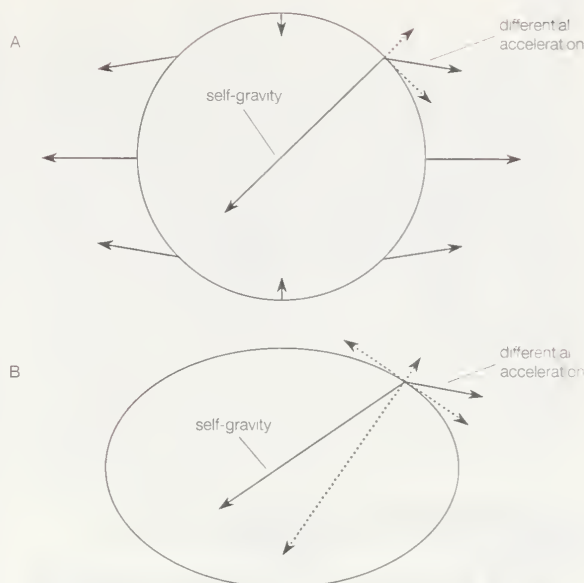


Figure 30: (A) Uncompensated tangential accelerations cause (B) tidal distortions in which all the differential accelerations are balanced by the change in the self-gravitational acceleration resulting from the distortion and by internal stresses (see text).

the differential acceleration due to the gravitational attraction of m_s , thus the elements do not fly off. If m_p were incompressible and perfectly rigid, the mass elements on the surface would weigh a little less than they would if m_s were not there but would not move relative to the centre of m_p . If m_p were fluid or otherwise not rigid, it would distort into an oval shape in the presence of m_s . The reason for this distortion is that the mass elements making up m_p that do not lie on the line joining the centres of m_p and m_s also experience a differential acceleration. Such differential accelerations are not perpendicular to the surface, however, and are therefore not compensated by the self-gravity that accelerates mass elements toward the centre of m_p . This is shown in Figure 30A, where one of the differential accelerations is resolved into two components (dotted arrows), one perpendicular and one tangential to the surface. The perpendicular component is compensated by the self-gravity; the tangential component is not. If m_p were entirely fluid, the uncompensated tangential components of the differential accelerations due to m_s would cause mass to flow toward the points on m_p that were either closest to m_s or farthest from m_s , until m_p would resemble Figure 30B. In this shape the self-gravity is no longer perpendicular to the surface but has a component opposite the tangential component of the differential acceleration. Only in this distorted shape will all the differential accelerations be compensated and the entire body accelerated like the centre. If m_p is not fluid but is rigid like rock or iron, part of the compensating acceleration will be provided by internal stress forces, and the body will distort less. As no material is perfectly rigid, there is always some tidal bulge, and compressibility of the material will further enhance this bulge. Note that the tidal distortion is independent of the orbital motion and would also occur if m_p and m_s were simply falling toward each other. (There is a similar tide raised on m_s by m_p that will be ignored for the present.)

If m_p rotates relative to m_s , an observer on the surface of m_p would successively rotate through the maxima and minima of the tidal distortion, which would tend to remain aligned with the direction to m_s . The observer would thereby experience two high and two low tides a day, as observed on Earth. Some of the energy of motion of any fluid parts of m_p and some of the energy stored as distortion of the solid parts as the tides wax and wane is converted into heat, and this dissipation of mechanical energy causes a delay in the response of the body to the tide-raising force. This means that high tide would occur at a given point on m_p as it rotates relative to m_s after m_s passes overhead. (On the Earth, the continents alter

High and low tides

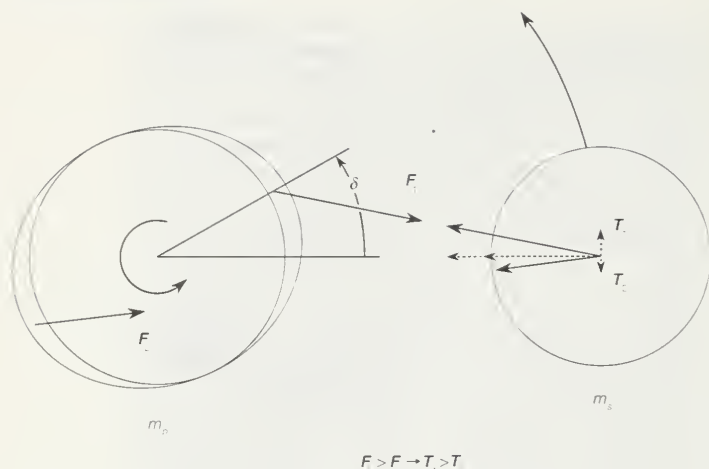


Figure 31: Unequal forces on two tidal bulges, leading to retardation of the spin of m_p and an acceleration of m_s in its orbit (see text).

the motion of the fluid ocean so much that ocean tides at continental coasts do not always lag the passage of the Moon overhead.) If m_p rotates in the same direction as m_s revolves in its orbit, the tidal bulge is carried ahead of m_s , as shown in Figure 31 by angle δ . Again, because the gravitational force between two masses varies as the inverse square of their separation, the tidal bulge closest to m_s experiences a greater attraction toward m_s (F_1 in Figure 31) than does the bulge farthest away (F_2). As these two forces are not aligned with the centre of m_p , there is a twisting effect, or torque, on m_p that retards its rate of rotation. This retardation will continue until the rotation is synchronous with the mean orbital motion of m_s . This has happened for the Moon, which keeps the same face toward the Earth.

From Newton's third law, there are equal and opposite forces acting on m_s corresponding to F_1 and F_2 . In Figure 31 these forces are represented as T_1 and T_2 , and each has been resolved into two components, one directed toward the centre of m_p and the other perpendicular to this direction. The inequality of these forces causes a net acceleration of m_s in its orbit, which thereby expands, as is observed for the Moon. Both the observed increase in the length of one day of 0.0016 second per century and the observed recession of the Moon of 3 to 4 centimetres per year are understood as consequences of the tides raised on the Earth.

In Figure 31, it has been assumed that the spin axis of m_p is perpendicular to the plane of the orbit of m_s . If the spin axis is inclined to this plane, the tidal bulge is carried out of the plane as well as ahead of m_s . This means that there is a twist, or torque, that changes the direction of the spin axis, so both the magnitude of spin and the direction of the spin axis experience a tidal evolution. The end point of tidal evolution for the spin state of one body of an isolated pair is rotation synchronous with the mean

orbital motion with the spin axis perpendicular to the orbit plane. This simple picture is complicated somewhat if other perturbations cause the orbital plane to precess. This precession for the lunar orbit causes its spin axis to be inclined $6^\circ 41'$ to the orbit plane as the end point of tidal evolution.

In addition to those of the Earth-Moon pair, numerous other consequences of tidal dissipation and the resulting evolution can be observed in the solar system and elsewhere in the Milky Way Galaxy. For example, all the major and close planetary satellites but one are observed to be rotating synchronously with their orbital motion. The exception is Saturn's satellite Hyperion. Tidal friction has indeed retarded Hyperion's initial spin rate to a value near that of synchronous rotation, but the combination of Hyperion's unusually asymmetric shape and its high orbital eccentricity leads to gravitational torques that make synchronous rotation unstable. As a result, the tides have brought Hyperion to a state where it tumbles chaotically with large changes in the direction and magnitude of its spin on time scales comparable to its orbital period of about 21 days.

Saturn-Hyperion tidal interactions

The assembly and maintenance of several orbital resonances among the satellites because of differential tidal expansion of the orbits have also been observed. The orbital resonances among Jupiter's satellites Io, Europa, and Ganymede, where the orbital periods are nearly in the ratio 1:2:4, maintain Io's orbital eccentricity at the value of 0.0041. This rather modest eccentricity causes sufficient variation in the magnitude and direction of Io's enormous tidal bulge to have melted a significant fraction of the satellite through dissipation of tidal energy in spite of Io's synchronous rotation. As a result, Io is the most volcanically active body in the solar system. The orbital eccentricity would normally be damped to zero by this large dissipation, but the orbital resonances with Europa and Ganymede prevent this from happening.

The distant planet Pluto and its satellite Charon are the only pair in the solar system that have almost certainly reached the ultimate end point where further tidal evolution has ceased altogether (the tiny tides raised by the Sun and other planets being neglected). In this state the orbit is circular, with both bodies rotating synchronously with the orbital motion and both spin axes perpendicular to the orbital plane.

The spin of the planet Mercury has been slowed by tides raised by the Sun to a final state where the spin angular velocity is exactly 1.5 times the orbital mean motion. This state is stable against further change because Mercury's high orbital eccentricity (0.206) allows restoring torques on the permanent (nontidal) axial asymmetry of the planet, which keeps the longest equatorial axis aligned with the direction to the Sun at perihelion. The tidal reduction of Mercury's average eccentricity (near 0.2) will cause insufficient change during the remaining lifetime of the Sun to disrupt this spin-orbit resonance. Finally, many close binary stars are observed to have circular orbits and synchronized spins—an example of tidal evolution elsewhere in the Milky Way Galaxy. (S.J.P.)

RELATIVISTIC MECHANICS

Relativistic mechanics is concerned with the motion of bodies whose relative velocities approach the speed of light c , or whose kinetic energies are comparable with the product of their masses m and the square of the velocity of light, or mc^2 . Such bodies are said to be relativistic, and when their motion is studied it is necessary to take into account Einstein's special theory of relativity. As long as gravitational effects can be ignored, which will be true so long as gravitational potential energy differences are small compared with mc^2 , the effects of Einstein's general theory of relativity may be safely ignored.

The bodies concerned may be sufficiently small that one may ignore their internal structure and size and regard them as point particles, in which case one speaks of relativistic point-particle mechanics; or one may need to take

into account their internal structure, in which case one speaks of relativistic continuum mechanics. This article is concerned only with relativistic point-particle mechanics. It is also assumed that quantum mechanical effects are unimportant, otherwise relativistic quantum mechanics or relativistic quantum field theory—the latter theory being a quantum mechanical extension of relativistic continuum mechanics—would have to be considered. The condition that allows quantum effects to be safely ignored is that the sizes and separations of the bodies concerned are larger than their Compton wavelengths. (The Compton wavelength of a body of mass m is given by h/mc , where h is Planck's constant.) Despite these restrictions, there are nevertheless a number of situations in nature where relativistic mechanics is applicable. For example, it is es-

sential to take into account the effects of relativity when calculating the motion of elementary particles accelerated to higher energies in particle accelerators, such as those at CERN (European Organization for Nuclear Research) near Geneva or at Fermilab (Fermi National Accelerator Laboratory) near Chicago. Moreover, such particles are caused to collide, thus creating further particles; although this creation process can only be understood through quantum mechanics, once the particles are well separated, they are subject to the laws of special relativity.

Similar remarks apply to cosmic rays that reach the Earth from outer space. In some cases, these have energies as high as 10^{20} electron volts (eV). An electron of that energy has a velocity that differs from that of light by about 1 part in 10^{28} , as can be seen from the relativistic relation between energy and velocity, which will be given later. For a proton of the same energy, the velocity would differ from that of light by about 1 part in 10^{22} . At a more mundane level, relativistic mechanics must be used to calculate the energies of electrons or positrons emitted by the decay of radioactive nuclei. Astrophysicists need to use relativistic mechanics when dealing with the energy sources of stars, the energy released in supernova explosions, and the motion of electrons moving in the atmospheres of pulsars or when considering the hot big bang. At temperatures in the very early universe above 10^{10} kelvins (K), at which typical thermal energies kT (where k is Boltzmann's constant and T is temperature) are comparable with the rest mass energy of the electron, the primordial plasma must have been relativistic. Relativistic mechanics also must be considered when dealing with satellite navigational systems used, for example, by the military, such as the Global Positioning System (GPS). In this case, however, it is the purely kinematic effect on the rate of clocks on board the satellites (*i.e.*, time dilation) that is important rather than the dynamic effects of relativity on the motion of the satellites themselves.

DEVELOPMENT OF THE SPECIAL THEORY OF RELATIVITY

Since the time of Galileo it has been realized that there exists a class of so-called inertial frames of reference—*i.e.*, in a state of uniform motion with respect to one another such that one cannot, by purely mechanical means, distinguish one from the other. It follows that the laws of mechanics must take the same form in every inertial frame of reference. To the accuracy of present-day technology, the class of inertial frames may be regarded as those that are neither accelerating nor rotating with respect to the distant galaxies. To specify the motion of a body relative to a frame of reference, one gives its position \mathbf{x} as a function of a time coordinate t (\mathbf{x} is called the position vector and has the components x , y , and z).

Newton's first law of motion (which remains true in special relativity) states that a body acted upon by no external forces will continue to move in a state of uniform motion relative to an inertial frame. It follows from this that the transformation between the coordinates (t, \mathbf{x}) and (t', \mathbf{x}') of two inertial frames with relative velocity \mathbf{u} must be related by a linear transformation. Before Einstein's special theory of relativity was published in 1905, it was usually assumed that the time coordinates measured in all inertial frames were identical and equal to an "absolute time." Thus,

$$t = t'. \tag{97}$$

The position coordinates \mathbf{x} and \mathbf{x}' were then assumed to be related by

$$\mathbf{x}' = \mathbf{x} - \mathbf{u}t. \tag{98}$$

The two formulas (97) and (98) are called a Galilean transformation. The laws of nonrelativistic mechanics take the same form in all frames related by Galilean transformations. This is the restricted, or Galilean, principle of relativity.

The position of a light-wave front speeding from the origin at time zero should satisfy

$$\mathbf{x}^2 - c^2t^2 = 0 \tag{99}$$

in the frame (t, \mathbf{x}) and

$$(\mathbf{x}')^2 - c^2(t')^2 = 0 \tag{100}$$

in the frame (t', \mathbf{x}') . Formula (100) does not transform into formula (99) using the Galilean transformations (97) and (98), however. Put another way, if one uses Galilean transformations one finds that the velocity of light depends on one's inertial frame, which is contrary to the Michelson-Morley experiment (see RELATIVITY). Einstein realized that either it is possible to determine a unique absolute frame of rest relative to which the motion of a light wave is given by equation (99) and its velocity is c only in that frame or the assumption that all inertial observers measure the same absolute time t —*i.e.*, formula (97)—must be wrong. Since he believed in (and experiment confirmed) the (extended) principle of relativity, which meant that one cannot, by any means, including the use of light waves, distinguish between two inertial frames in uniform relative motion, Einstein chose to give up the Galilean transformations (97) and (98) and replaced them with the Lorentz transformations:

Lorentz transformations

$$t' = \frac{(t - \mathbf{u} \cdot \mathbf{x}/c^2)}{\sqrt{1 - u^2/c^2}}, \tag{101}$$

$$\mathbf{x}'_{\parallel} = \frac{(\mathbf{x}_{\parallel} - \mathbf{u}t)}{\sqrt{1 - u^2/c^2}}. \tag{102a}$$

$$\mathbf{x}'_{\perp} = \mathbf{x}_{\perp}, \tag{102b}$$

where \mathbf{x}_{\parallel} and \mathbf{x}_{\perp} are the projections of \mathbf{x} parallel and perpendicular to the velocity \mathbf{u} , respectively, and similarly for \mathbf{x}' .

The reader may check that substitution of the Lorentz transformation formulas (101) and (102) into the left-hand side of equation (100) results in the left-hand side of equation (99). For simplicity, it has been assumed here and throughout this discussion, that the spatial axes are not rotated with respect to one another. Even in this case one sometimes considers Lorentz transformations that are more general than those of equations (101) and (102). These more general transformations may reverse the sense of time; *i.e.*, t and t' may have opposite signs or may reverse spatial orientation or parity. To distinguish this more general class of transformations from those of equations (101) and (102), one sometimes refers to (101) and (102) as proper Lorentz transformations.

The laws of light propagation are the same in all frames related by Lorentz transformations, and the velocity of light is the same in all such frames. The same is true of Maxwell's laws of electromagnetism. However, the usual laws of mechanics are not the same in all frames related by Lorentz transformations and thus must be altered to agree with the principle of relativity.

The unique absolute frame of rest with respect to which light waves had velocity c according to the prerelativistic viewpoint was often regarded, before Einstein, as being at rest relative to a hypothesized all-pervading ether. The vibrations of this ether were held to explain the phenomenon of electromagnetic radiation. The failure of experimenters to detect motion relative to this ether, together with the widespread acceptance of Einstein's special theory of relativity, led to the abandonment of the theory of the ether. It is ironic therefore to note that the discovery in 1964 by the American astrophysicists Arno Penzias and Robert Wilson of a universal cosmic microwave 3 K radiation background shows that the universe does indeed possess a privileged inertial frame. Nevertheless, this does not contradict special relativity because one cannot measure the Earth's velocity relative to it by experiments in a closed laboratory. One must actually detect the microwaves themselves.

If the relative velocity \mathbf{u} between inertial frames is small in magnitude compared with the velocity of light, then Galilean transformations and Lorentz transformations agree, as do the usual laws of nonrelativistic mechanics and the more accurate laws of relativistic mechanics. The requirement that the laws of physics take the same form in all inertial reference frames related by Lorentz transfor-

Inertial frames of reference

mations is called for the sake of brevity the requirement of relativistic invariance. It has become a powerful guide in the formation of new physical theories.

RELATIVISTIC SPACE-TIME

The modification of the usual laws of mechanics may be understood purely in terms of the Lorentz transformation formulas (101) and (102). It was pointed out, however, by the German mathematician Hermann Minkowski in 1908, that the Lorentz transformations have a simple geometric interpretation that is both beautiful and useful. The motion of a particle may be regarded as forming a curve made up of points, called events, in a four-dimensional space whose four coordinates comprise the three spatial coordinates $\mathbf{x} \equiv (x, y, z)$ and the time t .

Minkowski space-time

The four-dimensional space is called Minkowski space-time and the curve a world line. It is frequently useful to represent physical processes by space-time diagrams in which time runs vertically and the spatial coordinates run horizontally. Of course, since space-time is four-dimensional, at least one of the spatial dimensions in the diagram must be suppressed.

Newton's first law can be interpreted in four-dimensional space as the statement that the world lines of particles suffering no external forces are straight lines in space-time. Linear transformations take straight lines to straight lines, and Lorentz transformations have the additional property that they leave invariant the invariant interval τ through two events (t_1, \mathbf{x}_1) and (t_2, \mathbf{x}_2) given by

$$\tau^2 = (t_1 - t_2)^2 - \frac{(\mathbf{x}_1 - \mathbf{x}_2)^2}{c^2} \tag{103}$$

If the right-hand side of equation (103) is zero, the two events may be joined by a light ray and are said to be on each other's light cones because the light cone of any event (t, \mathbf{x}) in space-time is the set of points reachable from it by light rays (see Figure 32). Thus the set of all events (t_2, \mathbf{x}_2) satisfying equation (103) with zero on the right-hand side is the light cone of the event (t_1, \mathbf{x}_1) . Because Lorentz transformations leave invariant the space-time interval (103), all inertial observers agree on what the light cones are. In space-time diagrams it is customary to adopt a scaling of the time coordinate such that the light cones have a half angle of 45°.

If the right-hand side of equation (103) is strictly positive, in which case one says that the two events are timelike separated, or have a timelike interval, then one can find an

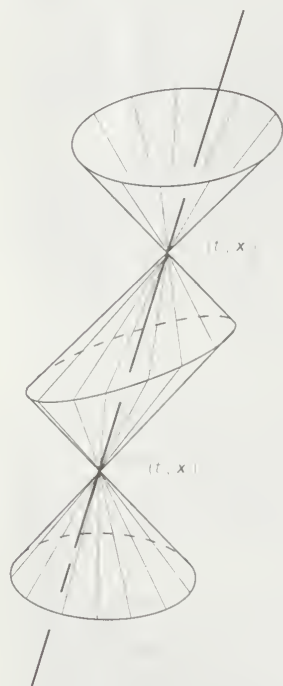


Figure 32: The world line of a particle traveling with speed less than that of light.

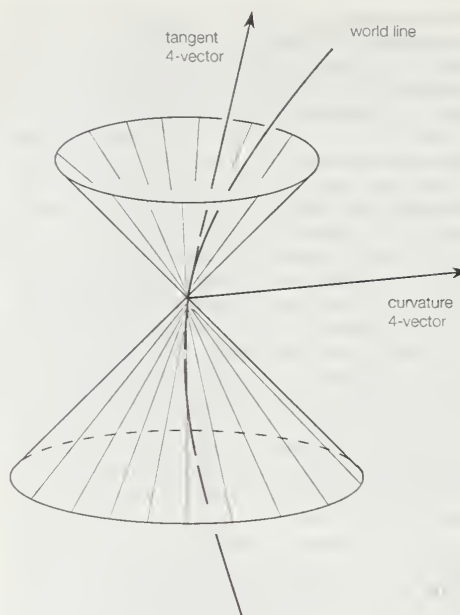


Figure 33: The world line of an accelerating body moving slower than the speed of light; the tangent vector corresponds to the body's 4-velocity and the curvature vector to its 4-acceleration.

inertial frame with respect to which the two events have the same spatial position. The straight world line joining the two events corresponds to the time axis of this inertial frame of reference. The quantity τ is equal to the difference in time between the two events in this inertial frame and is called the proper time between the two events. The proper time would be measured by any clock moving along the straight world line between the two events.

An accelerating body will have a curved world line that may be specified by giving its coordinates t and \mathbf{x} as a function of the proper time τ along the world line. The laws of either may be phrased in terms of the more familiar velocity $\mathbf{v} = d\mathbf{x}/dt$ and acceleration $\mathbf{a} = d^2\mathbf{x}/dt^2$ or in terms of the 4-velocity $(dt/d\tau, d\mathbf{x}/d\tau)$ and 4-acceleration $(d^2t/d\tau^2, d^2\mathbf{x}/d\tau^2)$. Just as an ordinary vector like \mathbf{v} has three components, $v_x, v_y,$ and $v_z,$ a 4-vector has four components. Geometrically the 4-velocity and 4-acceleration correspond, respectively, to the tangent vector and the curvature vector of the world line (see Figure 33). If the particle moves slower than light, the tangent, or velocity, vector at each event on the world line points inside the light cone of that event, and the acceleration, or curvature, vector points outside the light cone. If the particle moves with the speed of light, then the tangent vector lies on the light cone at each event on the world line. The proper time τ along a world line moving with a speed less than light is not an independent quantity from t and \mathbf{x} : it satisfies

$$\left(\frac{dt}{d\tau}\right)^2 - \frac{1}{c^2}\left(\frac{d\mathbf{x}}{d\tau}\right)^2 = 1. \tag{104}$$

For a particle moving with exactly the speed of light, one cannot define a proper time τ . One can, however, define a so-called affine parameter that satisfies equation (104) with zero on the right-hand side. For the time being this discussion will be restricted to particles moving with speeds less than light.

Equation (104) does not fix the sign of τ relative to that of t . It is usual to resolve this ambiguity by demanding that the proper time τ increase as the time t increases. This requirement is invariant under Lorentz transformations of the form of equations (101) and (102). The tangent vector then points inside the future light cone and is said to be future-directed and timelike (see Figure 34). One may if one wishes attach an arrow to the world line to indicate this fact. One says that the particle moves forward in time. It was pointed out by the Swiss physicist Ernest C.G. Stückelberg de Breidenbach and by the American physicist Richard Feynman that a meaning can be attached to world lines moving backward in time—i.e., for those for

4-vectors

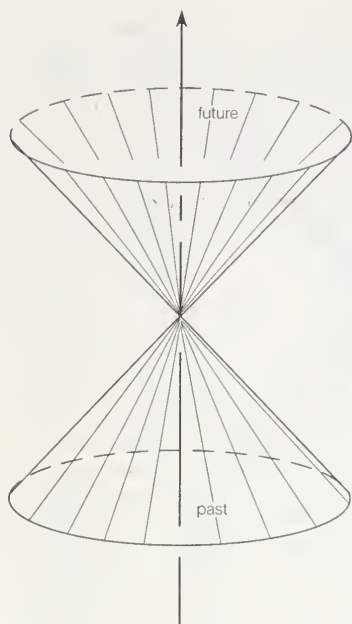


Figure 34: The world line of a particle moving forward in time (see text).

which ordinary time t decreases as proper time τ increases. Since, as shall be shown later, the energy E of a particle is $mc^2 dt/d\tau$, such world lines correspond to the motion of particles with negative energy. It is possible to interpret these world lines in terms of antiparticles, as will be seen when particles moving in a background electromagnetic field are considered.

The fundamental laws of motion for a body of mass m in relativistic mechanics are

$$m \frac{d^2 t}{d\tau^2} = f^0 \tag{105}$$

and

$$m \frac{d^2 \mathbf{x}}{d\tau^2} = \mathbf{f} \tag{106}$$

where m is the constant so-called rest mass of the body and the quantities (f^0, \mathbf{f}) are the components of the force 4-vector. Equations (105) and (106), which relate the curvature of the world line to the applied forces, are the same in all inertial frames related by Lorentz transformations. The quantities $(m dt/d\tau, m d\mathbf{x}/d\tau)$ make up the 4-momentum of the particle. According to Minkowski's reformulation of special relativity, a Lorentz transformation may be thought of as a generalized rotation of points of Minkowski space-time into themselves. It induces an identical rotation on the 4-acceleration and force 4-vectors. To say that both of these 4-vectors experience the same generalized rotation or Lorentz transformation is simply to say that the fundamental laws of motion (105) and (106) are the same in all inertial frames related by Lorentz transformations. Minkowski's geometric ideas provided a powerful tool for checking the mathematical consistency of special relativity and for calculating its experimental consequences. They also have a natural generalization in the general theory of relativity, which incorporates the effects of gravity.

RELATIVISTIC MOMENTUM, MASS, AND ENERGY

The law of motion (106) may also be expressed as:

$$\frac{d}{dt} \left(\frac{m\mathbf{v}}{\sqrt{1 - v^2/c^2}} \right) = \mathbf{F}, \tag{107}$$

where $\mathbf{F} = \mathbf{f} \sqrt{1 - v^2/c^2}$. Equation (107) is of the same form as Newton's second law of motion, which states that the rate of change of momentum equals the applied force. \mathbf{F} is the Newtonian force, but the Newtonian relation between momentum \mathbf{p} and velocity \mathbf{v} in which $\mathbf{p} = m\mathbf{v}$ is modified to become

$$\mathbf{p} = \frac{m\mathbf{v}}{\sqrt{1 - v^2/c^2}}. \tag{108}$$

Consider a relativistic particle with positive energy and electric charge q moving in an electric field \mathbf{E} and magnetic field \mathbf{B} ; it will experience an electromagnetic, or Lorentz, force given by $\mathbf{F} = q\mathbf{E} + q\mathbf{v} \times \mathbf{B}$. If $t(\tau)$ and $\mathbf{x}(\tau)$ are the time and space coordinates of the particle, it follows from equations (105) and (106), with $f^0 = (q\mathbf{E} \cdot \mathbf{v})dt/d\tau$ and $\mathbf{f} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B})dt/d\tau$, that $-t(-\tau)$ and $-\mathbf{x}(-\tau)$ are the coordinates of a particle with positive energy and the opposite electric charge $-q$ moving in the same electric and magnetic field. A particle of the opposite charge but with the same rest mass as the original particle is called the original particle's antiparticle. It is in this sense that Feynman and Stückelberg spoke of antiparticles as particles moving backward in time. This idea is a consequence of special relativity alone. It really comes into its own, however, when one considers relativistic quantum mechanics.

Anti-particles

Just as in nonrelativistic mechanics, the rate of work done when the point of application of a force \mathbf{F} is moved with velocity \mathbf{v} equals $\mathbf{F} \cdot \mathbf{v}$ when measured with respect to the time coordinate t . This work goes into increasing the energy E of the particle. Taking the dot product of equation (107) with \mathbf{v} gives

$$\frac{dE}{dt} = \mathbf{F} \cdot \mathbf{v}, \tag{109a}$$

where

$$E = \frac{mc^2}{\sqrt{1 - v^2/c^2}} = mc^2 \frac{dt}{d\tau}. \tag{109b}$$

The reader should note that the 4-momentum is just $(E/c^2, \mathbf{p})$. It was once fairly common to encounter the use of a "velocity-dependent mass" equal to E/c^2 . However, experience has shown that its introduction serves no useful purpose and may lead to confusion, and it is not used in this article. The invariant quantity is the rest mass m . For that reason it has not been thought necessary to add a subscript or superscript to m to emphasize that it is the rest mass rather than a velocity-dependent quantity. When subscripts are attached to a mass, they indicate the particular particle of which it is the rest mass.

If the applied force \mathbf{F} is perpendicular to the velocity \mathbf{v} , it follows from equation (109) that the energy E , or, equivalently, the velocity squared v^2 , will be constant, just as in Newtonian mechanics. This will be true, for example, for a particle moving in a purely magnetic field with no electric field present. It then follows from equation (107) that the shape of the orbits of the particle are the same according to the classical and the relativistic equations. However, the rate at which the orbits are traversed differs according to the two theories. If w is the speed according to the nonrelativistic theory and v that according to special relativity, then $w = v \sqrt{1 - v^2/c^2}$.

For velocities that are small compared with that of light,

$$E \approx mc^2 + \frac{1}{2} m v^2. \tag{110}$$

The first term, mc^2 , which remains even when the particle is at rest, is called the rest mass energy. For a single particle, its inclusion in the expression for energy might seem to be a matter of convention: it appears as an arbitrary constant of integration. However, for systems of particles that undergo collisions, its inclusion is essential.

Both theory and experiment agree that, in a process in which particles of rest masses m_1, m_2, \dots, m_n collide or decay or transmute one into another, both the total energy $E_1 + E_2 + \dots + E_n$ and the total momentum $\mathbf{p}_1 + \mathbf{p}_2 + \dots + \mathbf{p}_n$ are the same before and after the process, even though the number of particles may not be the same before and after. This corresponds to conservation of the total 4-momentum $(E_1 + E_2 + \dots + E_n)/c^2, \mathbf{p}_1 + \mathbf{p}_2 + \dots + \mathbf{p}_n$.

The relativistic law of energy-momentum conservation thus combines and generalizes in one relativistically invariant expression the separate conservation laws of prerelativistic physics: the conservation of mass, the conservation

Law of energy-momentum conservation

Rotation in space-time

of momentum, and the conservation of energy. In fact, the law of conservation of mass becomes incorporated in the law of conservation of energy and is modified if the amount of energy exchanged is comparable with the rest mass energy of any of the particles.

For example, if a particle of mass M at rest decays into two particles the sum of whose rest masses $m_1 + m_2$ is smaller than M (see Figure 35), then the two momenta p_1 and p_2 must be equal in magnitude and opposite in direction. The quantity $T = E - mc^2$ is the kinetic energy of the particle. In such a decay the initial kinetic energy is zero. Since the conservation of energy implies that in the process $Mc^2 = T_1 + T_2 + m_1c^2 + m_2c^2$, one speaks of the conversion of an amount $(M - m_1 - m_2)c^2$ of rest mass energy to kinetic energy. It is precisely this process that provides the large amount of energy available during nuclear fission, for example, in the spontaneous fission of the uranium-235 isotope. The opposite process occurs in nuclear fusion when two particles fuse to form a particle of smaller total rest mass. The difference $(m_1 + m_2 - M)$ multiplied by c^2 is called the binding energy. If the two initial particles are both at rest, a fourth particle is required to satisfy the conservation of energy and momentum. The rest mass of this fourth particle will not change, but it will acquire kinetic energy equal to the binding energy minus the kinetic energy of the fused particles. Perhaps the most important examples are the conversion of hydrogen to helium in the centre of stars, such as the Sun, and during thermonuclear reactions used in atomic bombs.

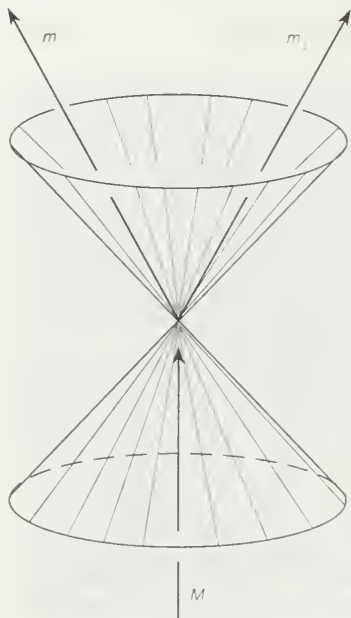


Figure 35: The decay of a particle of mass M into two particles the sum of whose rest masses is less than M (see text).

This article has so far dealt only with particles with non-vanishing rest mass whose velocities must always be less than that of light. One may always find an inertial reference frame with respect to which they are at rest and their energy in that frame equals mc^2 . However, special relativity allows a generalization of classical ideas to include particles with vanishing rest masses that can move only with

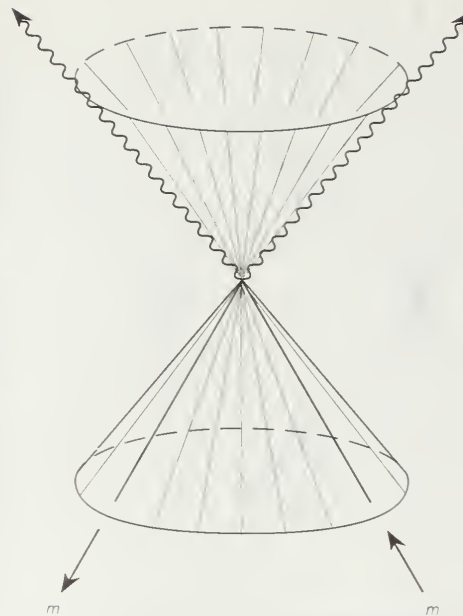


Figure 36: The world lines of an electron (moving forward in time) and a positron (moving backward in time) that annihilate into two photons (see text).

the velocity of light. Particles in nature that correspond to this possibility and that could not, therefore, be incorporated into the classical scheme are the photon, which is associated with the transmission of electromagnetic radiation, three species of neutrinos associated with the weak interaction responsible for radioactive decay, and—more speculatively—the graviton, which plays the same role with respect to gravitational waves as does the photon with respect to electromagnetic waves. Strictly speaking, it is not yet known whether the rest masses of the three neutrino species vanish, but in many processes their energies exceed their possible rest masses by so much that they may be regarded as effectively massless. The velocity v of any particle in relativistic mechanics is given by $v = pc^2/E$, and the relation between energy E and momentum is $E^2 = m^2c^4 + p^2c^2$. Thus for massless particles $E = |p|c$ and the 4-momentum is given by $(|p|/c, p)$. It follows from the relativistic laws of energy and momentum conservation that, if a massless particle were to decay, it could do so only if the particles produced were all strictly massless and their momenta p_1, p_2, \dots, p_n were all strictly aligned with the momentum p of the original massless particle. Since this is a situation of vanishing likelihood, it follows that strictly massless particles are absolutely stable.

It also follows that one or more massive particles cannot decay into a single massless particle, conserving both energy and momentum. They can, however, decay into two or more massless particles, and indeed this is observed in the decay of the neutral pion into photons and in the annihilation of an electron and a positron pair into photons. In the latter case, the world lines of the annihilating particles meet at the space-time event where they annihilate. Using the interpretation of Feynman and Stückelberg, one may view these two world lines as a single continuous world line with two portions, one moving forward in time and one moving backward in time (see Figure 36). This interpretation plays an important role in the quantum theory of such processes. (G.W.G.)

Stability
of
massless
particles

MECHANICS OF SOLIDS

The application of the principles of mechanics to bulk matter is conventionally divided into the mechanics of fluids and the mechanics of solids. The entire subject is often called continuum mechanics, particularly when matter is assumed to be continuously divisible and no reference is made to its discrete structure at microscopic scales well below those of the application or phenomenon of interest.

Solid mechanics is concerned with the stressing, deformation, and failure of solid materials and structures. What, then, is a solid? Any material, fluid or solid, can support normal forces. These are forces directed perpendicular, or normal, to a material plane across which they act. The force per unit of area of that plane is called the normal stress. Water at the base of a pond, air in an automobile

Definition
of a solid

ture, the stones of a Roman arch, rocks at the base of a mountain, the skin of a pressurized airplane cabin, a stretched rubber band, and the bones of a runner all support force in that way (some only when the force is compressive). A material is called solid rather than fluid if it can also support a substantial shearing force over the time scale of some natural process or technological application of interest. Shearing forces are directed parallel, rather than perpendicular, to the material surface on which they act; the force per unit of area is called shear stress. For example, consider a vertical metal rod that is fixed to a support at its upper end and has a weight attached at its lower end. If one considers a horizontal surface through the material of the rod, it will be evident that the rod supports normal stress. But it also supports shear stress, and this becomes evident when one considers the forces carried across a plane that is neither horizontal nor vertical through the rod. Thus, while water and air provide no long-term support of shear stress, granite, steel, and rubber normally do so and are therefore called solids. Materials with tightly bound atoms or molecules, such as the crystals formed below melting temperature by most substances or simple compounds and the amorphous structures formed in glass and many polymer substances at sufficiently low temperature, are usually considered solids.

The distinction between solids and fluids is not precise and in many cases will depend on the time scale. Consider the hot rocks of the Earth's mantle. When a large earthquake occurs, an associated deformation disturbance called a seismic wave propagates through the adjacent rock, and the entire Earth is set into vibrations which, following a sufficiently large earthquake, may remain detectable with precise instruments for several weeks. The rocks of the mantle are then described as solid—as they would also be on the time scale of, say, tens to thousands of years, over which stresses rebuild enough in the source region to cause one or a few repetitions of the earthquake. But on a significantly longer time scale, say, on the order of a million years, the hot rocks of the mantle are unable to support shearing stresses and flow as a fluid. The substance called Silly Putty (trademark), a polymerized silicone gel familiar to many children, is another example. If a ball of it is left to sit on a table at room temperature, it flows and flattens on a time scale of a few minutes to an hour. But if picked up and tossed as a ball against a wall, so that large forces act only over the short time of the impact, the Silly Putty bounces back and retains its shape like a highly elastic solid.

Types of
solid deformation

Several types of solids can be distinguished according to their mechanical behaviour. In the simple but common case when a solid material is loaded at a sufficiently low temperature or short time scale, and with sufficiently limited stress magnitude, its deformation is fully recovered upon unloading. The material is then said to be elastic. But substances can also deform permanently, so that not all the deformation is recovered. For example, if one bends a metal coat hanger substantially and then releases the loading, it springs back only partially toward its initial shape; it does not fully recover but remains bent. The metal of the coat hanger has been permanently deformed, and in this case, for which the permanent deformation is not so much a consequence of longtime loading at sufficiently high temperature but more a consequence of subjecting the material to large stresses (above the yield stress), the permanent deformation is described as a plastic deformation and the material is called elastic-plastic. Permanent deformation of a sort that depends mainly on time of exposure to a stress—and that tends to increase significantly with time of exposure—is called viscous, or creep, deformation, and materials that exhibit those characteristics, as well as tendencies for elastic response, are called viscoelastic solids (or sometimes viscoplastic solids, when the permanent strain is emphasized rather than the tendency for partial recovery of strain upon unloading).

Solid mechanics has many applications. All those who seek to understand natural phenomena involving the stressing, deformation, flow, and fracture of solids, as well as all those who would have knowledge of such phenomena to improve living conditions and accomplish human

objectives, have use for solid mechanics. The latter activities are, of course, the domain of engineering, and many important modern subfields of solid mechanics have been actively developed by engineering scientists concerned, for example, with mechanical, structural, materials, civil, or aerospace engineering. Natural phenomena involving solid mechanics are studied in geology, seismology, and tectonophysics, in materials science and the physics of condensed matter, and in some branches of biology and physiology. Furthermore, because solid mechanics poses challenging mathematical and computational problems, it (as well as fluid mechanics) has long been an important topic for applied mathematicians concerned, for example, with partial differential equations and with numerical techniques for digital computer formulations of physical problems.

Here is a sampling of some of the issues addressed using solid mechanics concepts: How do flows develop in the Earth's mantle and cause continents to move and ocean floors to subduct (*i.e.*, be thrust) slowly beneath them? How do mountains form? What processes take place along a fault during an earthquake, and how do the resulting disturbances propagate through the Earth as seismic waves, shaking, and perhaps collapsing, buildings and bridges? How do landslides occur? How does a structure on a clay soil settle with time, and what is the maximum bearing pressure that the footing of a building can exert on a soil or rock foundation without rupturing it? What materials should be chosen, and how should their proportion, shape, and loading be controlled, to make safe, reliable, durable, and economical structures—whether airframes, bridges, ships, buildings, chairs, artificial heart valves, or computer chips—and to make machinery such as jet engines, pumps, and bicycles? How do vehicles (cars, planes, ships) respond by vibration to the irregularity of surfaces or mediums along which they move, and how are vibrations controlled for comfort, noise reduction, and safety against fatigue failure? How rapidly does a crack grow in a cyclically loaded structure, whether a bridge, engine, or airplane wing or fuselage, and when will it propagate catastrophically? How can the deformability of structures during impact be controlled so as to design crashworthiness into vehicles? How are the materials and products of a technological civilization formed—*e.g.*, by extruding metals or polymers through dies, rolling material into sheets, punching out complex shapes, and so on? By what microscopic processes do plastic and creep strains occur in polycrystals? How can different materials, such as fibre-reinforced composites, be fashioned together to achieve combinations of stiffness and strength needed in specific applications? What is the combination of material properties and overall response needed in downhill skis or in a tennis racket? How does the human skull respond to impact in an accident? How do heart muscles control the pumping of blood in the human body, and what goes wrong when an aneurysm develops?

Issues
addressed
by solid
mechanics

History

Solid mechanics developed in the outpouring of mathematical and physical studies following the great achievement of Newton in stating the laws of motion, although it has earlier roots. The need to understand and control the fracture of solids seems to have been a first motivation. Leonardo da Vinci sketched in his notebooks a possible test of the tensile strength of a wire. Galileo, who died in the year of Newton's birth (1642), had investigated the breaking loads of rods under tension and concluded that the load was independent of length and proportional to the cross section area, this being a first step toward a concept of stress. He also investigated the breaking loads on beams that were suspended horizontally from a wall into which they were built.

Concepts of stress, strain, and elasticity. The English scientist Robert Hooke discovered in 1660, but published only in 1678, that for many materials the displacement under a load was proportional to force, thus establishing the notion of (linear) elasticity but not yet in a way that was expressible in terms of stress and strain. Edme Mariotte in France published similar discoveries in 1680 and,

Relation
between
stress and
strain

in addition, reached an understanding of how beams like those studied by Galileo resist transverse loadings—or, more precisely, resist the torques caused by those transverse loadings—by developing extensional and compressional deformations, respectively, in material fibres along their upper and lower portions. It was for the Swiss mathematician and mechanician Jakob Bernoulli to observe, in the final paper of his life, in 1705, that the proper way of describing deformation was to give force per unit area, or stress, as a function of the elongation per unit length, or strain, of a material fibre under tension. The Swiss mathematician and mechanician Leonhard Euler, who was taught mathematics by Jakob's brother Johann Bernoulli, proposed, among many contributions, a linear relation between stress σ and strain ϵ , in 1727, of the form $\sigma = E\epsilon$, where the coefficient E is now generally called Young's modulus after the British naturalist Thomas Young, who developed a related idea in 1807.

The notion that there is an internal tension acting across surfaces in a deformed solid was expressed by the German mathematician and physicist Gottfried Wilhelm Leibniz in 1684 and Jakob Bernoulli in 1691. Also, Jakob Bernoulli and Euler introduced the idea that at a given section along the length of a beam there were internal tensions amounting to a net force and a net torque (see below). Euler introduced the idea of compressive normal stress as the pressure in a fluid in 1752. The French engineer and physicist Charles-Augustin Coulomb was apparently the first to relate the theory of a beam as a bent elastic line to stress and strain in an actual beam, in a way never quite achieved by Bernoulli and, although possibly recognized, never published by Euler. He developed the famous expression $\sigma = My/I$ for the stress due to the pure bending of a homogenous linear elastic beam; here M is the torque, or bending moment, y is the distance of a point from an axis that passes through the section centroid, parallel to the torque axis, and I is the integral of y^2 over the section area. The French mathematician Antoine Parent introduced the concept of shear stress in 1713, but Coulomb was the one who extensively developed the idea, first in connection with beams and with the stressing and failure of soil in 1773 and then in studies of frictional slip in 1779.

Contributions of
Cauchy

It was the great French mathematician Augustin-Louis Cauchy, originally educated as an engineer, who in 1822 formalized the concept of stress in the context of a generalized three-dimensional theory, showed its properties as consisting of a 3×3 symmetric array of numbers that transform as a tensor, derived the equations of motion for a continuum in terms of the components of stress, and developed the theory of linear elastic response for isotropic solids. As part of his work in this area, Cauchy also introduced the equations that express the six components of strain (three extensional and three shear) in terms of derivatives of displacements for the case in which all those derivatives are much smaller than unity; similar expressions had been given earlier by Euler in expressing rates of straining in terms of the derivatives of the velocity field in a fluid.

Beams, columns, plates, and shells. The 1700s and early 1800s were a productive period during which the mechanics of simple elastic structural elements were developed—well before the beginnings in the 1820s of the general three-dimensional theory. The development of beam theory by Euler, who generally modeled beams as elastic lines that resist bending, as well as by several members of the Bernoulli family and by Coulomb, remains among the most immediately useful aspects of solid mechanics, in part for its simplicity and in part because of the pervasiveness of beams and columns in structural technology. Jakob Bernoulli proposed in his final paper of 1705 that the curvature of a beam was proportional to its bending moment. Euler in 1744 and Johann's son, Daniel Bernoulli, in 1751 used the theory to address the transverse vibrations of beams, and in 1757 Euler gave his famous analysis of the buckling of an initially straight beam subjected to a compressive loading; such a beam is commonly called a column. Following a suggestion of Daniel Bernoulli in 1742, Euler in 1744 introduced the

concept of strain energy per unit length for a beam and showed that it is proportional to the square of the beam's curvature. Euler regarded the total strain energy as the quantity analogous to the potential energy of a discrete mechanical system. By adopting procedures that were becoming familiar in analytical mechanics and following from the principle of virtual work as introduced in 1717 by Johann Bernoulli for such discrete systems as pin-connected rigid bodies, Euler rendered the energy stationary and in this way developed the calculus of variations as an approach to the equations of equilibrium and motion of elastic structures.

That same variational approach played a major role in the development by French mathematicians in the early 1800s of a theory of small transverse displacements and vibrations of elastic plates. This theory was developed in preliminary form by Sophie Germain and was also worked on by Siméon-Denis Poisson in the early 1810s; they considered a flat plate as an elastic plane that resists curvature. Claude-Louis-Marie Navier gave a definitive development of the correct energy expression and governing differential equation a few years later. An uncertainty of some duration arose in the theory from the fact that the final partial differential equation for the transverse displacement is such that it is impossible to prescribe, simultaneously, along an unsupported edge of the plate, both the twisting moment per unit length of middle surface and the transverse shear force per unit length. This was finally resolved in 1850 by the Prussian physicist Gustav Robert Kirchhoff, who applied virtual work and variational calculus procedures in the framework of simplifying kinematic assumptions that fibres initially perpendicular to the plate's middle surface remain so after deformation of that surface.

The first steps in the theory of thin shells were taken by Euler in the 1770s; he addressed the deformation of an initially curved beam as an elastic line and provided a simplified analysis of the vibration of an elastic bell as an array of annular beams. Johann's grandson, Jakob Bernoulli "the Younger," further developed this model in the last year of his life as a two-dimensional network of elastic lines, but he could not develop an acceptable treatment. Shell theory did not attract attention again until a century after Euler's work. The first consideration of shells from a three-dimensional elastic viewpoint was advanced by Hermann Aron in 1873. Acceptable thin-shell theories for general situations, appropriate for cases of small deformation, were then developed by the British mathematician, mechanician, and geophysicist Augustus Edward Hough Love in 1888 and by the British mathematician and physicist Horace Lamb in 1890 (there is no uniquely correct theory, as the Dutch applied mechanician and engineer W.T. Koiter and the Soviet mechanician V.V. Novozhilov clarified in the 1950s; the difference between predictions of acceptable theories is small when the ratio of shell thickness to a typical length scale is small). Shell theory remained of immense interest well beyond the mid-1900s, in part because so many problems lay beyond the linear theory (rather small transverse displacements often dramatically alter the way that a shell supports load by a combination of bending and membrane action) and in part because of the interest in such lightweight structural forms for aeronautical technology.

The general theory of elasticity. Linear elasticity as a general three-dimensional theory began to be developed in the early 1820s based on Cauchy's work. Simultaneously, Navier had developed an elasticity theory based on a simple corpuscular, or particle, model of matter in which particles interacted with their neighbours by a central force attraction between particle pairs. As was gradually realized, following work by Navier, Cauchy, and Poisson in the 1820s and '30s, the particle model is too simple and makes predictions concerning relations among elastic moduli that are not met by experiment. Most of the subsequent development of this subject was in terms of the continuum theory. Controversies concerning the maximum possible number of independent elastic moduli in the most general anisotropic solid were settled by the British mathematician George Green in 1837. Green pointed out that the existence of an elastic strain energy

required that of the 36 elastic constants relating the 6 stress components to the 6 strains, at most 21 could be independent. The Scottish physicist Lord Kelvin put this consideration on sounder ground in 1855 as part of his development of macroscopic thermodynamics, showing that a strain energy function must exist for reversible isothermal or adiabatic (isentropic) response and working out such results as the (very modest) temperature changes associated with isentropic elastic deformation (see below *Thermodynamic considerations*).

Applications of elasticity

The middle and late 1800s were a period in which many basic elastic solutions were derived and applied to technology and to the explanation of natural phenomena. The French mathematician Adhémar-Jean-Claude Barré de Saint-Venant derived in the 1850s solutions for the torsion of noncircular cylinders, which explained the necessity of warping displacement of the cross section in the direction parallel to the axis of twisting, and for the flexure of beams due to transverse loadings; the latter allowed understanding of approximations inherent in the simple beam theory of Jakob Bernoulli, Euler, and Coulomb. The German physicist Heinrich Rudolf Hertz developed solutions for the deformation of elastic solids as they are brought into contact and applied these to model details of impact collisions. Solutions for stress and displacement due to concentrated forces acting at an interior point of a full space were derived by Kelvin, and those on the surface of a half space by the French mathematician Joseph Valentin Boussinesq and the Italian mathematician Valentino Cerruti. The Prussian mathematician Leo August Pochhammer analyzed the vibrations of an elastic cylinder, and Lamb and the Prussian physicist Paul Jaerisch derived the equations of general vibration of an elastic sphere in the 1880s, an effort that was continued by many seismologists in the 1900s to describe the vibrations of the Earth. In 1863 Kelvin had derived the basic form of the solution of the static elasticity equations for a spherical solid, and these were applied in following years to such problems as calculating the deformation of the Earth due to rotation and tidal forcing and measuring the effects of elastic deformability on the motions of the Earth's rotation axis.

The classical development of elasticity never fully confronted the problem of finite elastic straining, in which material fibres change their lengths by other than very small amounts. Possibly this was because the common materials of construction would remain elastic only for very small strains before exhibiting either plastic straining or brittle failure. However, natural polymeric materials show elasticity over a far wider range (usually also with enough time or rate effects that they would more accurately be characterized as viscoelastic), and the widespread use of natural rubber and similar materials motivated the development of finite elasticity. While many roots of the subject were laid in the classical theory, especially in the work of Green, Gabrio Piola, and Kirchhoff in the mid-1800s, the development of a viable theory with forms of stress-strain relations for specific rubbery elastic materials, as well as an understanding of the physical effects of the nonlinearity in simple problems such as torsion and bending, was mainly the achievement of the British-born engineer and applied mathematician Ronald S. Rivlin in the 1940s and '50s.

Waves. Poisson, Cauchy, and George G. Stokes showed that the equations of the general theory of elasticity predicted the existence of two types of elastic deformation waves which could propagate through isotropic elastic solids. These are called body waves. In the faster type, called longitudinal, dilational, or irrotational waves, the particle motion is in the same direction as that of wave propagation; in the slower type, called transverse, shear, or rotational waves, it is perpendicular to the propagation direction. No analogue of the shear wave exists for propagation through a fluid medium, and that fact led seismologists in the early 1900s to understand that the Earth has a liquid core (at the centre of which there is a solid inner core).

Lord Rayleigh showed in 1885 that there is a wave type that could propagate along surfaces, such that the mo-

tion associated with the wave decayed exponentially with distance into the material from the surface. This type of surface wave, now called a Rayleigh wave, propagates typically at slightly more than 90 percent of the shear wave speed and involves an elliptical path of particle motion that lies in planes parallel to that defined by the normal to the surface and the propagation direction. Another type of surface wave, with motion transverse to the propagation direction and parallel to the surface, was found by Love for solids in which a surface layer of material sits atop an elastically stiffer bulk solid; this defines the situation for the Earth's crust. The shaking in an earthquake is communicated first to distant places by body waves, but these spread out in three dimensions and to conserve the energy propagated by the wave field must diminish in their displacement amplitudes as r^{-1} , where r is the distance from the source. The surface waves spread out in only two dimensions and must, for the same reason, diminish only as fast as $r^{-1/2}$. Thus, the shaking effect of the surface waves from a crustal earthquake is normally felt more strongly, and is potentially more damaging, at moderate to large distances. Indeed, well before the theory of waves in solids was in hand, Thomas Young had suggested in his 1807 lectures on natural philosophy that the shaking of an earthquake "is probably propagated through the earth in the same manner as noise is conveyed through air." (It had been suggested by the American mathematician and astronomer John Winthrop, following his experience of the "Boston" earthquake of 1755, that the ground shaking was due to a disturbance propagated like sound through the air.)

Types of surface waves

With the development of ultrasonic transducers operated on piezoelectric principles, the measurement of the reflection and scattering of elastic waves has developed into an effective engineering technique for the nondestructive evaluation of materials for detection of such potentially dangerous defects as cracks. Also, very strong impacts, whether from meteorite collision, weaponry, or blasting and the like in technological endeavours, induce waves in which material response can be well outside the range of linear elasticity, involving any or all of finite elastic strain, plastic or viscoplastic response, and phase transformation. These are called shock waves; they can propagate much beyond the speed of linear elastic waves and are accompanied by significant heating.

Stress concentrations and fracture. In 1898 G. Kirsch derived the solution for the stress distribution around a circular hole in a much larger plate under remotely uniform tensile stress. The same solution can be adapted to the tunnelling cylindrical cavity of a circular section in a bulk solid. Kirsch's solution showed a significant concentration of stress at the boundary, by a factor of three when the remote stress was uniaxial tension. Then in 1907 the Russian mathematician Gury Vasilyevich Kolosov, and independently in 1914 the British engineer Charles Edward Inglis, derived the analogous solution for stresses around an elliptical hole. Their solution showed that the concentration of stress could become far greater, as the radius of curvature at an end of the hole becomes small compared with the overall length of the hole. These results provided the insight to sensitize engineers to the possibility of dangerous stress concentrations at sharp reentrant corners, notches, cutouts, keyways, screw threads, and similar openings in structures for which the nominal stresses were at otherwise safe levels. Such stress concentration sites are places from which a crack can nucleate.

The elliptical hole of Kolosov and Inglis defines a crack in the limit when one semimajor axis goes to zero, and the Inglis solution was adopted by the British aeronautical engineer A.A. Griffith in 1921 to describe a crack in a brittle solid. In that work Griffith made his famous proposition that a spontaneous crack growth would occur when the energy released from the elastic field just balanced the work required to separate surfaces in the solid. Following a hesitant beginning, in which Griffith's work was initially regarded as important only for very brittle solids such as glass, there developed, largely under the impetus of the American engineer and physicist George R. Irwin, a major body of work on the mechanics of crack growth and

Crack growth

fracture, including fracture by fatigue and stress corrosion cracking, starting in the late 1940s and continuing into the 1990s. This was driven initially by the cracking of a number of American Liberty ships during World War II, by the failures of the British Comet airplane, and by a host of reliability and safety issues arising in aerospace and nuclear reactor technology. The new complexion of the subject extended beyond the Griffith energy theory and, in its simplest and most widely employed version in engineering practice, used Irwin's stress intensity factor as the basis for predicting crack growth response under service loadings in terms of laboratory data that is correlated in terms of that factor. That stress intensity factor is the coefficient of a characteristic singularity in the linear elastic solution for the stress field near a crack tip; it is recognized as providing a proper characterization of crack tip stressing in many cases, even though the linear elastic solution must be wrong in detail near the crack tip owing to nonelastic material response, large strain, and discreteness of material microstructure.

Dislocations. The Italian elastician and mathematician Vito Volterra introduced in 1905 the theory of the elastostatic stress and displacement fields created by dislocating solids. This involves making a cut in a solid, displacing its surfaces relative to one another by some fixed amount, and joining the sides of the cut back together, filling in with material as necessary. The initial status of this work was simply regarded as an interesting way of generating elastic fields, but, in the early 1930s, Geoffrey Ingram Taylor, Egon Orowan, and Michael Polanyi realized that just such a process could be going on in ductile crystals and could provide an explanation of the low plastic shear strength of typical ductile solids, much as Griffith's cracks explained low fracture strength under tension. In this case, the displacement on the dislocated surface corresponds to one atomic lattice spacing in the crystal. It quickly became clear that this concept provided the correct microscopic description of metal plasticity, and, starting with Taylor in the 1930s and continuing into the 1990s, the use of solid mechanics to explore dislocation interactions and the microscopic basis of plastic flow in crystalline materials has been a major topic, with many distinguished contributors.

The mathematical techniques advanced by Volterra are now in common use by earth scientists in explaining ground displacement and deformation induced by tectonic faulting. Also, the first elastodynamic solutions for the rapid motion of crystal dislocations, developed by South African materials scientist F.R.N. Nabarro in the early 1950s, were quickly adapted by seismologists to explain the radiation from propagating slip distributions on faults. The Japanese seismologist H. Nakano had already shown in 1923 how to represent the distant waves radiated by an earthquake as the elastodynamic response to a pair of force dipoles amounting to zero net torque. (All his manuscripts were destroyed in the fire in Tokyo associated with the great Kwanto earthquake in that same year, but copies of some had been sent to Western colleagues and the work survived.)

Continuum plasticity theory. The macroscopic theory of plastic flow has a history nearly as old as that of elasticity. While in the microscopic theory of materials, the word "plasticity" is usually interpreted as denoting deformation by dislocation processes, in macroscopic continuum mechanics it is taken to denote any type of permanent deformation of materials, especially those of a type for which time or rate of deformation effects are not the most dominant feature of the phenomenon (the terms viscoplasticity, creep, or viscoelasticity are usually used in such cases). Coulomb's work of 1773 on the frictional yielding of soils under shear and normal stress has been mentioned; yielding denotes the occurrence of large shear deformations without significant increase in applied stress. His results were used to explain the pressure of soils against retaining walls and footings in the work of the French mathematician and engineer Jean Victor Poncelet in 1840 and the Scottish engineer and physicist William John Macquorn Rankine in 1853. The inelastic deformation of soils and rocks often takes place in situations for which the deforming mass is infiltrated by groundwater,

and Austrian-American civil engineer Karl Terzaghi in the 1920s developed the concept of effective stress, whereby the stresses that enter a criterion of yielding or failure are not the total stresses applied to the saturated soil or rock mass but rather the effective stresses, which are the difference between the total stresses and those of a purely hydrostatic stress state with pressure equal to that in the pore fluid. Terzaghi also introduced the concept of consolidation, in which the compression of a fluid-saturated soil can take place only as the fluid slowly flows through the pore space under pressure gradients, according to Darcy's law; this effect accounts for the time-dependent settlement of constructions over clay soils.

Apart from the earlier observation of plastic flow at large stresses in the tensile testing of bars, the theory of continuum plasticity for metallic materials begins with Henri Edouard Tresca in 1864. His experiments on the compression and indentation of metals led him to propose that this type of plasticity, in contrast to that in soils, was essentially independent of the average normal stress in the material and dependent only on shear stresses, a feature later rationalized by the dislocation mechanism. Tresca proposed a yield criterion for macroscopically isotropic metal polycrystals based on the maximum shear stress in the material, and that was used by Saint-Venant to solve an early elastic-plastic problem, that of the partly plastic cylinder in torsion, and also to solve for the stresses in a completely plastic tube under pressure.

The German applied mechanician Ludwig Prandtl developed the rudiments of the theory of plane plastic flow in 1920 and 1921, with an analysis of indentation of a ductile solid by a flat-ended rigid indenter, and the resulting theory of plastic slip lines was completed by H. Hencky in 1923 and Hilda Geiringer in 1930. Additional developments include the methods of plastic limit analysis, which allowed engineers to directly calculate upper and lower bounds to the plastic collapse loads of structures or to forces required in metal forming. Those methods developed gradually over the early 1900s on a largely intuitive basis, first for simple beam structures and later for plates, and were put on a rigorous basis within the rapidly developing mathematical theory of plasticity about 1950 by Daniel C. Drucker and William Prager in the United States and Rodney Hill in Great Britain.

The Austrian-American applied mathematician Richard von Mises proposed in 1913 that a mathematically simpler theory of plasticity than that based on the Tresca yield criterion could be based on the second tensor invariant of the deviatoric stresses (*i.e.*, of the total stresses minus those of a hydrostatic state in which pressure is equal to the average normal stress over all planes). An equivalent yield criterion had been proposed independently by the Polish engineer Maksymilian Tytus Huber. The Mises theory incorporates a proposal by M. Levy in 1871 that components of the plastic strain increment tensor are in proportion to one another just as are the components of deviatoric stress. This criterion was generally found to provide slightly better agreement with experiment than did that of Tresca, and most work on the application of plasticity theory uses this form. Following a suggestion of Prandtl, E. Reuss completed the theory in 1930 by adding an elastic component of strain increments, related to stress increments in the same way as for linear elastic response. This formulation was soon generalized to include strain hardening, whereby the value of the second invariant for continued yielding increases with ongoing plastic deformation, and was extended to high-temperature creep response in metals or other hot solids by assuming that the second invariant of the plastic (now generally called "creep") strain rate is a function of that same invariant of the deviatoric stress, typically a power law type with Arrhenius temperature dependence.

This formulation of plastic and viscoplastic, or creep, response has been applied to all manner of problems in materials and structural technology and in flow of geologic masses. Representative problems addressed include the growth and subsequent coalescence of microscopic voids in the ductile fracture of metals, the theory of the indentation hardness test, the extrusion of metal rods and

Microscopic description of plastic flow

Plane plastic flow

Strain hardening

rolling of metal sheets, design against collapse of ductile steel structures, estimation of the thickness of the Greenland Ice Sheet, and modeling the geologic evolution of the Plateau of Tibet. Other types of elastic-plastic theories intended for analysis of ductile single crystals originate from the work of G.I. Taylor and Hill and base the criterion for yielding on E. Schmid's concept from the 1920s of a critical resolved shear stress along a crystal slip plane, in the direction of an allowed slip on that plane; this sort of yield condition has approximate support from the dislocation theory of plasticity.

Viscoelasticity. The German physicist Wilhelm Weber noticed in 1835 that a load applied to a silk thread produced not only an immediate extension but also a continuing elongation of the thread with time. This type of viscoelastic response is especially notable in polymeric solids but is present to some extent in all types of solids and often does not have a clear separation from what could be called viscoplastic, or creep, response. In general, if all of the strain is ultimately recovered when a load is removed from a body, the response is termed viscoelastic, but the term is also used in cases for which sustained loading leads to strains that are not fully recovered. The Austrian physicist Ludwig Boltzmann developed in 1874 the theory of linear viscoelastic stress-strain relations. In their most general form, these involve the notion that a step loading (a suddenly imposed stress that is subsequently maintained constant) causes an immediate strain followed by a time-dependent strain which, for different materials, either may have a finite limit at long time or may increase indefinitely with time. Within the assumption of linearity, the strain at time t in response to a general time-dependent stress history $\sigma(t')$ can then be written as the sum (or integral) of terms that involve the step-loading strain response due to a step loading $dt' d\sigma(t')/dt'$ at time t' . The theory of viscoelasticity is important for consideration of the attenuation of stress waves and the damping of vibrations.

A new class of problems arose with the mechanics of very-long-molecule polymers, which do not have significant cross-linking and exist either in solution or as a melt. These are fluids in the sense that they cannot long support shear stress, but at the same time they have remarkable properties like those of finitely deformed elastic solids. A famous demonstration is to pour one of these fluids slowly from a beaker and to cut the flowing stream suddenly with scissors; if the cut is not too far below the place of exit from the beaker, the stream of falling fluid immediately contracts elastically and returns to the beaker. The molecules are elongated during flow but tend to return to their thermodynamically preferred coiled configuration when forces are removed.

The theory of such materials came under intense development in the 1950s after the British applied mathematician James Gardner Oldroyd showed in 1950 how viscoelastic stress-strain relations of a memory type could be generalized to a flowing fluid. This requires that the constitutive relation, or rheological relation, between the stress history and the deformation history at a material "point" be properly invariant to a superposed history of rigid rotation, which should not affect the local physics determining that relation (the resulting Coriolis and centrifugal effects are quite negligible at the scale of molecular interactions). Important contributions on this issue were made by the applied mathematicians Stanislaw Zaremba and Gustav Andreas Johannes Jaumann in the first decade of the 1900s; they showed how to make tensorial definitions of stress rate that were invariant to superposed spin and thus were suitable for use in constitutive relations. But it was only during the 1950s that these concepts found their way into the theory of constitutive relations for general viscoelastic materials; independently, a few years later, properly invariant stress rates were adopted in continuum formulations of elastic-plastic response.

Computational mechanics. The digital computer revolutionized the practice of many areas of engineering and science, and solid mechanics was among the first fields to benefit from its impact. Many computational techniques have been used in this field, but the one that emerged

by the end of 1970s as, by far, the most widely adopted is the finite-element method. This method was outlined by the mathematician Richard Courant in 1943 and was developed independently, and put to practical use on computers, in the mid-1950s by the aeronautical structures engineers M.J. Turner, Ray W. Clough, Harold Clifford Martin, and LeRoy J. Topp in the United States and J.H. Argyris and Sydney Kelsey in Britain. Their work grew out of earlier attempts at systematic structural analysis for complex frameworks of beam elements. The method was soon recast in a variational framework and related to earlier efforts at deriving approximate solutions of problems described by variational principles. The new technique involved substituting trial functions of unknown amplitude into the variational functional, which is then rendered stationary as an algebraic function of the amplitude coefficients. In the most common version of the finite-element method, the domain to be analyzed is divided into cells, or elements, and the displacement field within each element is interpolated in terms of displacements at a few points around the element boundary (and sometimes within it) called nodes. The interpolation is done so that the displacement field is continuous across element boundaries for any choice of the nodal displacements. The strain at every point can thus be expressed in terms of nodal displacements, and it is then required that the stresses associated with these strains, through the stress-strain relations of the material, satisfy the principle of virtual work for arbitrary variation of the nodal displacements. This generates as many simultaneous equations as there are degrees of freedom in the finite element model, and numerical techniques for solving such systems of equations are programmed for computer solution.

Basic principles

In addressing any problem in continuum or solid mechanics, three factors must be considered: (1) the Newtonian equations of motion, in the more general form recognized by Euler, expressing conservation of linear and angular momentum for finite bodies (rather than just for point particles), and the related concept of stress, as formalized by Cauchy, (2) the geometry of deformation and thus the expression of strains in terms of gradients in the displacement field, and (3) the relations between stress and strain that are characteristic of the material in question, as well as of the stress level, temperature, and time scale of the problem considered.

These three considerations suffice for most problems. They must be supplemented, however, for solids undergoing diffusion processes in which one material constituent moves relative to another (which may be the case for fluid-infiltrated soils or petroleum reservoir rocks) and in cases for which the induction of a temperature field by deformation processes and the related heat transfer cannot be neglected. These cases require that the following also be considered: (4) equations for conservation of mass of diffusing constituents, (5) the first law of thermodynamics, which introduces the concept of heat flux and relates changes in energy to work and heat supply, and (6) relations that express the diffusive fluxes and heat flow in terms of spatial gradients of appropriate chemical potentials and of temperature. In many important technological devices, electric and magnetic fields affect the stressing, deformation, and motion of matter. Examples are provided by piezoelectric crystals and other ceramics for electric or magnetic actuators and by the coils and supporting structures of powerful electromagnets. In these cases, two more considerations must be added: (7) James Clerk Maxwell's set of equations interrelating electric and magnetic fields to polarization and magnetization of material media and to the density and motion of electric charge, and (8) augmented relations between stress and strain, which now, for example, express all of stress, polarization, and magnetization in terms of strain, electric field, magnetic intensity, and temperature. The second law of thermodynamics, combined with the above-mentioned principles, serves to constrain physically allowed relations between stress, strain, and temperature in (3) and also

Considerations in solid mechanics problems

Long-molecule polymers

constrains the other types of relations described in (6) and (8) above. Such expressions, which give the relationships between stress, deformation, and other variables, are commonly referred to as constitutive relations.

In general, the stress-strain relations are to be determined by experiment. A variety of mechanical testing machines and geometric configurations of material specimens have been devised to measure them. These allow, in different cases, simple tensile, compressive, or shear stressing, and sometimes combined stressing with several different components of stress, as well as the determination of material response over a range of temperatures, strain rates, and loading histories. The testing of round bars under tensile stress, with precise measurement of their extension to obtain the strain, is common for metals and for technological ceramics and polymers. For rocks and soils, which generally carry load in compression, the most common test involves a round cylinder that is compressed along its axis, often while being subjected to confining pressure on its curved face. Frequently, a measurement interpreted by solid mechanics theory is used to determine some of the properties entering stress-strain relations. For example, measuring the speed of deformation waves or the natural frequencies of vibration of structures can be used to extract the elastic moduli of materials of known mass density, and measurement of indentation hardness of a metal can be used to estimate its plastic shear strength.

In some favourable cases, stress-strain relations can be calculated approximately by applying principles of mechanics at the microscale of the material considered. In a composite material, the microscale could be regarded as the scale of the separate materials making up the reinforcing fibres and matrix. When their individual stress-strain relations are known from experiment, continuum mechanics principles applied at the scale of the individual constituents can be used to predict the overall stress-strain relations for the composite. For rubbery polymer materials, made up of long chain molecules that randomly configure themselves into coil-like shapes, some aspects of the elastic stress-strain response can be obtained by applying principles of statistical thermodynamics to the partial uncoiling of the array of molecules by imposed strain. For a single crystallite of an element such as silicon or aluminum or for a simple compound like silicon carbide, the relevant microscale is that of the atomic spacing in the crystals; quantum mechanical principles governing atomic force laws at that scale can be used to estimate elastic constants. In the case of plastic flow processes in metals and in sufficiently hot ceramics, the relevant microscale involves the network of dislocation lines that move within crystals. These lines shift atom positions relative to one another by one atomic spacing as they move along slip planes. Important features of elastic-plastic and viscoplastic stress-strain relations can be understood by modeling the stress dependence of dislocation generation and motion and the resulting dislocation entanglement and immobilization processes that account for strain hardening.

To examine the mathematical structure of the theory, considerations (1) to (3) above will now be further developed. For this purpose, a continuum model of matter will be used, with no detailed reference to its discrete structure at molecular—or possibly other larger microscopic—scales far below those of the intended application.

LINEAR AND ANGULAR MOMENTUM PRINCIPLES: STRESS AND EQUATIONS OF MOTION

Let \mathbf{x} denote the position vector of a point in space as measured relative to the origin of a Newtonian reference frame; \mathbf{x} has the components (x_1, x_2, x_3) relative to a Cartesian set of axes, which is fixed in the reference frame and denoted as the 1, 2, and 3 axes in Figure 37. Suppose that a material occupies the part of space considered, and let $\mathbf{v} = \mathbf{v}(\mathbf{x}, t)$ be the velocity vector of the material point that occupies position \mathbf{x} at time t ; that same material point will be at position $\mathbf{x} + \mathbf{v}dt$ in an infinitesimal interval dt later. Let $\rho = \rho(\mathbf{x}, t)$ be the mass density of the material. Here \mathbf{v} and ρ are macroscopic variables. What is idealized in the continuum model as a material point, moving as a smooth function of time, will correspond on molecular-

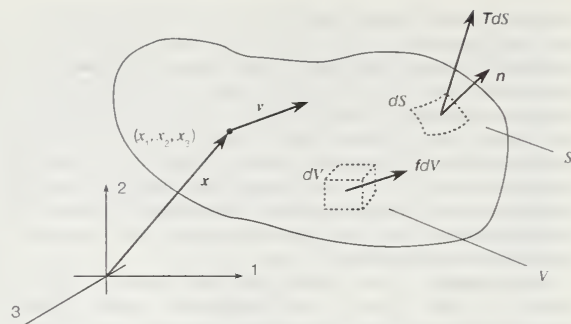


Figure 37: The position vector \mathbf{x} and the velocity vector \mathbf{v} of a material point, the body force $f dV$ acting on an element dV of volume, and the surface force $T dS$ acting on an element dS of surface in a Cartesian coordinate system 1, 2, 3 (see text).

length (or larger but still “microscopic”) scales to a region with strong fluctuations of density and velocity. In terms of phenomena at such scales, ρ corresponds to an average of mass per unit of volume, and $\rho\mathbf{v}$ to an average of linear momentum per unit volume, as taken over spatial and temporal scales that are large compared to those of the microscale processes but still small compared to those of the intended application or phenomenon under study. Thus, from the microscopic viewpoint, \mathbf{v} of the continuum theory is a mass-weighted average velocity.

The linear momentum \mathbf{P} and angular momentum \mathbf{H} (relative to the coordinate origin) of the matter instantaneously occupying any volume V of space are then given by summing up the linear and angular momentum vectors of each element of material. Such summation over infinitesimal elements is represented mathematically by the integrals $\mathbf{P} = \int_V \rho \mathbf{v} dV$ and $\mathbf{H} = \int_V \rho \mathbf{x} \times \mathbf{v} dV$. In this discussion attention is limited to situations in which relativistic effects can be ignored. Let \mathbf{F} denote the total force and \mathbf{M} the total torque, or moment (relative to the coordinate origin), acting instantaneously on the material occupying any arbitrary volume V . The basic laws of Newtonian mechanics are the linear and angular momentum principles that $\mathbf{F} = d\mathbf{P}/dt$ and $\mathbf{M} = d\mathbf{H}/dt$, where time derivatives of \mathbf{P} and \mathbf{H} are calculated following the motion of the matter that occupies V at time t . When either \mathbf{F} or \mathbf{M} vanishes, these equations of motion correspond to conservation of linear or angular momentum.

An important, very common, and nontrivial class of problems in solid mechanics involves determining the deformed and stressed configuration of solids or structures that are in static equilibrium; in that case the relevant basic equations are $\mathbf{F} = 0$ and $\mathbf{M} = 0$. The understanding of such conditions for equilibrium, at least in a rudimentary form, long predates Newton. Indeed, Archimedes of Syracuse (3rd century BC), the great Greek mathematician and arguably the first theoretically and experimentally minded physical scientist, understood these equations at least in a nonvectorial form appropriate for systems of parallel forces. This is shown by his treatment of the hydrostatic equilibrium of a partially submerged body and by his establishment of the principle of the lever (torques about the fulcrum sum to zero) and the concept of centre of gravity.

Stress. Assume that \mathbf{F} and \mathbf{M} derive from two types of forces, namely, body forces \mathbf{f} , such as gravitational attractions—defined such that force $f dV$ acts on volume element dV (see Figure 37)—and surface forces, which represent the mechanical effect of matter immediately adjoining that along the surface S of the volume V being considered. Cauchy formalized in 1822 a basic assumption of continuum mechanics that such surface forces could be represented as a stress vector \mathbf{T} , defined so that $T dS$ is an element of force acting over the area dS of the surface (Figure 37). Hence, the principles of linear and angular momentum take the forms

$$\int_S T dS + \int_V f dV = \mathbf{F} = \frac{d\mathbf{P}}{dt} = \int_V \rho a dV; \quad \text{and} \quad (111)$$

$$\int_S \mathbf{x} \times T dS + \int_V \mathbf{x} \times f dV = \mathbf{M} = \frac{d\mathbf{H}}{dt} = \int_V \rho \mathbf{x} \times a dV. \quad (112)$$

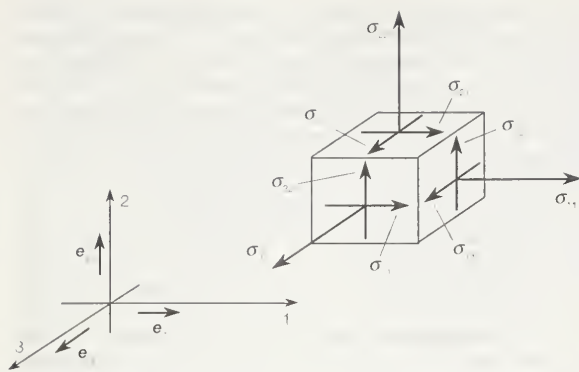


Figure 38: The nine components of a stress tensor. The first index denotes the direction of the normal, or perpendicular, stresses to the plane across which the contact force acts, and the second index denotes the direction of the component of force (see text).

which are now assumed to hold good for every conceivable choice of region V . In calculating the right-hand sides, which come from dP/dt and dH/dt , it has been noted that ρdV is an element of mass and is therefore time-invariant; also, $\mathbf{a} = \mathbf{a}(\mathbf{x}, t) = d\mathbf{v}/dt$ is the acceleration, where the time derivative of \mathbf{v} is taken following the motion of a material point so that $\mathbf{a}(\mathbf{x}, t)dt$ corresponds to the difference between $\mathbf{v}(\mathbf{x} + \mathbf{v}dt, t + dt)$ and $\mathbf{v}(\mathbf{x}, t)$. A more detailed analysis of this step shows that the understanding of what TdS denotes must now be adjusted to include averages, over temporal and spatial scales that are large compared to those of microscale fluctuations, of transfers of momentum across the surface S due to the microscopic fluctuations about the motion described by the macroscopic velocity \mathbf{v} .

The nine quantities σ_{ij} ($i, j = 1, 2, 3$) are called stress components; these will vary with position and time—i.e., $\sigma_{ij} = \sigma_{ij}(\mathbf{x}, t)$ —and have the following interpretation. Consider an element of surface dS through a point \mathbf{x} with dS oriented so that its outer normal (pointing away from the region V , bounded by S) points in the positive x_i direction, where i is any of 1, 2, or 3. Then σ_{i1} , σ_{i2} , and σ_{i3} at \mathbf{x} are defined as the Cartesian components of the stress vector \mathbf{T} (called $\mathbf{T}^{(i)}$) acting on this dS . Figure 38 shows the components of such stress vectors for faces in each of the three coordinate directions. To use a vector notation with \mathbf{e}_1 , \mathbf{e}_2 , and \mathbf{e}_3 denoting unit vectors along the coordinate axes (Figure 38), $\mathbf{T}^{(i)} = \sigma_{i1}\mathbf{e}_1 + \sigma_{i2}\mathbf{e}_2 + \sigma_{i3}\mathbf{e}_3$. Thus, the stress σ_{ij} at \mathbf{x} is the stress in the j direction associated with an i -oriented face through point \mathbf{x} ; the physical dimension of the σ_{ij} is [force]/[length]². The components σ_{11} , σ_{22} , and σ_{33} are stresses directed perpendicular, or normal, to the face on which they act and are normal stresses; the σ_{ij} with $i \neq j$ are directed parallel to the face on which they act and are shear stresses.

By hypothesis, the linear momentum principle applies for

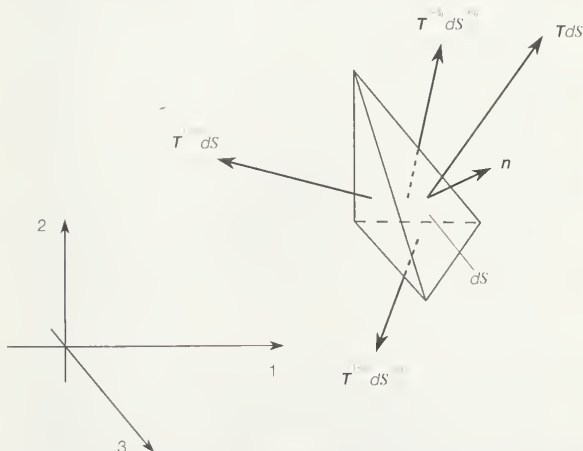


Figure 39: The force TdS acting on an arbitrarily inclined face (whose outward unit normal vector is \mathbf{n}). Stress vectors $\mathbf{T}^{(-1)}$, $\mathbf{T}^{(-2)}$, and $\mathbf{T}^{(-3)}$ act on the faces perpendicular to the coordinate axes.

any volume V . Consider a small tetrahedron (Figure 39) at \mathbf{x} with an inclined face having an outward unit normal vector \mathbf{n} and its other three faces oriented perpendicular to the three coordinate axes. Letting the size of the tetrahedron shrink to zero, the linear momentum principle requires that the stress vector \mathbf{T} on a surface element with outward normal \mathbf{n} be expressed as a linear function of the σ_{ij} at \mathbf{x} . The relation is such that the j component of the stress vector \mathbf{T} is $T_j = n_1\sigma_{1j} + n_2\sigma_{2j} + n_3\sigma_{3j}$, for ($j = 1, 2, 3$). This relation for \mathbf{T} (or T_j) also demonstrates that the σ_{ij} have the mathematical property of being the components of a second-rank tensor.

Second-rank tensor

Suppose that a different set of Cartesian reference axes $1'$, $2'$, and $3'$ have been chosen. Let x'_1 , x'_2 , and x'_3 denote the components of the position vector of point \mathbf{x} and let σ'_{kl} ($k, l = 1, 2, 3$) denote the nine stress components relative to that coordinate system. The σ'_{kl} can be written as the 3×3 matrix $[\sigma']$, and the σ_{ij} as the matrix $[\sigma]$, where the first index is the matrix row number and the second is the column number. Then the expression for T_j implies that $[\sigma'] = [\alpha][\sigma][\alpha]^T$, which is the defining equation of a second-rank tensor. Here $[\alpha]$ is the orthogonal transformation matrix, having components $\alpha_{pq} = \mathbf{e}'_p \cdot \mathbf{e}_q$, $q = 1, 2, 3$ and satisfying $[\alpha]^T[\alpha] = [\alpha][\alpha]^T = [I]$, where the superscript T denotes transpose (interchange rows and columns) and $[I]$ denotes the unit matrix, a 3×3 matrix with unity for every diagonal element and zero elsewhere; also, the matrix multiplications are such that if $[A] = [B][C]$, then $A_{ij} = B_{ik}C_{kj} + B_{2k}C_{2j} + B_{3k}C_{3j}$.

Equations of motion. Now the linear momentum principle may be applied to an arbitrary finite body. Using the expression for T_j above and the divergence theorem of multivariable calculus, which states that integrals over the area of a closed surface S , with integrand $n_i f(\mathbf{x})$, may be rewritten as integrals over the volume V enclosed by S , with integrand $\partial f(\mathbf{x})/\partial x_i$; when $f(\mathbf{x})$ is a differentiable function, one may derive that

$$\frac{\partial \sigma_{1j}}{\partial x_1} + \frac{\partial \sigma_{2j}}{\partial x_2} + \frac{\partial \sigma_{3j}}{\partial x_3} + f_j = \rho a_j \quad (j = 1, 2, 3), \quad (113)$$

at least when the σ_{ij} are continuous and differentiable, which is the typical case. These are the equations of motion for a continuum. Once the above consequences of the linear momentum principle are accepted, the only further result that can be derived from the angular momentum principle is that $\sigma_{ij} = \sigma_{ji}$ ($i, j = 1, 2, 3$). Thus, the stress tensor is symmetric.

Principal stresses. Symmetry of the stress tensor has the important consequence that, at each point \mathbf{x} , there exist three mutually perpendicular directions along which there are no shear stresses. These directions are called the principal stress directions, and the corresponding normal stresses are called the principal stresses. If the principal stresses are ordered algebraically as σ_1 , σ_{II} , and σ_{III} (Figure 40), then the normal stress on any face (given as $\sigma_n = \mathbf{n} \cdot \mathbf{T}$) satisfies $\sigma_1 \leq \sigma_n \leq \sigma_{III}$. The principal stresses are the eigenvalues (or characteristic values) s , and the principal directions the eigenvectors \mathbf{n} , of the problem $\mathbf{T} = s\mathbf{n}$, or $[\sigma]\{\mathbf{n}\} = s\{\mathbf{n}\}$ in matrix notation with the 3-column $\{\mathbf{n}\}$ representing \mathbf{n} . It has solutions when $\det([\sigma] - s[I]) = -s^3 + I_1 s^2 + I_2 s + I_3 = 0$, with $I_1 = \text{tr}[\sigma]$, $I_2 = -(1/2)I_1^2 + (1/2)\text{tr}([\sigma][\sigma])$, and $I_3 = \det[\sigma]$. Here "det" denotes determinant and "tr" denotes trace, or sum of diagonal elements, of a matrix. Since the principal stresses are determined by I_1 , I_2 , and I_3 and can have no dependence on how one chooses the coordinate system with respect to which the components of stress are referred, I_1 , I_2 , and I_3 must be independent of that choice and are therefore called stress invariants. One may readily verify that they have the same values when evaluated in terms of σ'_{ij} above as in terms of σ_{ij} by using the tensor transformation law and properties noted for the orthogonal transformation matrix.

Stress invariants

Very often, in both nature and technology, there is interest in structural elements in forms that might be identified as strings, wires, rods, bars, beams, or columns, or as membranes, plates, or shells. These are usually idealized as, respectively, one- or two-dimensional continua. One possible approach is then to develop the consequences

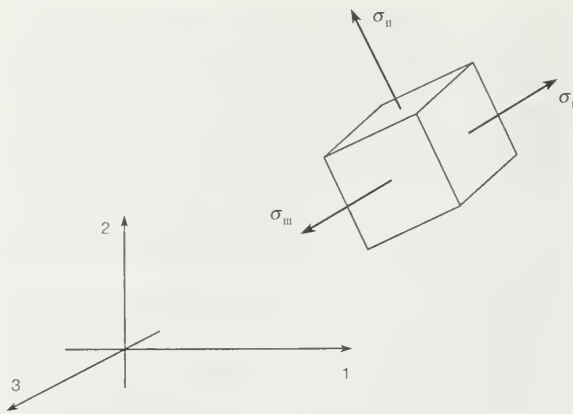


Figure 40: Principal stresses (see text).

of the linear and angular momentum principles entirely within that idealization, working in terms of net axial and shear forces and bending and twisting torques at each point along a one-dimensional continuum, or in terms of forces and torques per unit length of surface in a two-dimensional continuum.

GEOMETRY OF DEFORMATION

Strain and strain-displacement relations. The shape of a solid or structure changes with time during a deformation process. To characterize deformation, or strain, a certain reference configuration is adopted and called undeformed. Often, that reference configuration is chosen as an unstressed state, but such is neither necessary nor always convenient. If time is measured from zero at a moment when the body exists in that reference configuration, then the upper case X may be used to denote the position vectors of material points when $t = 0$. At some other time t , a material point that was at X will have moved to some spatial position x . The deformation is thus described as the mapping $x = x(X, t)$, with $x = x(X, 0) = X$. The displacement vector u is then $u = x(X, t) - X$; also, $v = \partial x(X, t) / \partial t$ and $a = \partial^2 x(X, t) / \partial t^2$.

It is simplest to write equations for strain in a form that, while approximate in general, is suitable for the case when any infinitesimal line element dX of the reference configuration undergoes extremely small rotations and fractional change in length, in deforming to the corresponding line element dx . These conditions are met when $|\partial u_i / \partial X_j| \ll 1$. Many solids are often sufficiently rigid, at least under the loadings typically applied to them, that these conditions are realized in practice. Linearized expressions for strain in terms of $[\partial u / \partial X]$, appropriate to this situation, are called small strain or infinitesimal strain. Expressions for strain will also be given that are valid for rotations and fractional length changes of arbitrary magnitude; such expressions are called finite strain.

Two simple types of strain are extensional strain and shear strain. Consider a rectangular parallelepiped, a brick-like block of material with mutually perpendicular planar faces, and let the edges of the block be parallel to the 1, 2, and 3 axes. If the block is deformed homogeneously, so that each planar face moves perpendicular to itself and so that the faces remain orthogonal (*i.e.*, the parallelepiped is deformed into another rectangular parallelepiped), then the block is said to have undergone extensional strain relative to these axes. If the edge lengths of the undeformed parallelepiped are denoted as ΔX_1 , ΔX_2 , and ΔX_3 , and those of the deformed parallelepiped as Δx_1 , Δx_2 , and Δx_3 (see Figure 41A, where the dashed-line figure represents the reference configuration and the solid-line figure the deformed configuration), then the quantities $\lambda_1 = \Delta x_1 / \Delta X_1$, $\lambda_2 = \Delta x_2 / \Delta X_2$, and $\lambda = \Delta x_3 / \Delta X_3$ are called stretch ratios. There are various ways that extensional strain can be defined in terms of them. Note that the change in displacement in, say, the x_1 direction between points at one end of the block and those at the other is $\Delta u_1 = (\lambda_1 - 1)\Delta X_1$. For example, if E_{11} denotes the extensional strain along the x_1 direction, then the most commonly understood defini-

tion of strain is $E_{11} = (\text{change in length}) / (\text{initial length}) = (\Delta x_1 - \Delta X_1) / \Delta X_1 = \Delta u_1 / \Delta X_1 = \lambda_1 - 1$. A variety of other measures of extensional strain can be defined by $E_{11} = g(\lambda_1)$, where the function $g(\lambda)$ satisfies $g(1) = 0$ and $g'(1) = 1$, so as to agree with the above definition when λ_1 is very near 1. Two such measures in common use are the strain $E_{11}^m = (\lambda_1^2 - 1) / 2$, based on the change of metric tensor, and the logarithmic strain $E_{11}^l = \ln(\lambda_1)$.

To define a simple shear strain, consider the same rectangular parallelepiped, but now deform it so that every point on a plane of type $X_2 = \text{constant}$ moves only in the x_1 direction by an amount that increases linearly with X_2 . Thus, the deformation $x_1 = \gamma X_2 + X_1$, $x_2 = X_2$, $x_3 = X_3$ defines a homogeneous simple shear strain of amount γ and is illustrated in Figure 41B. Note that this strain causes no change of volume. For small strain, the shear strain γ can be identified as the reduction in angle between two initially perpendicular lines.

Simple shear strain

Small-strain tensor. The small strains, or infinitesimal strains, ϵ_{ij} are appropriate for situations with $|\partial u_k / \partial X_l| \ll 1$ for all k and l . Two infinitesimal material fibres, one initially in the 1 direction and the other in the 2 direction, are shown in Figure 42 as dashed lines in the reference configuration and as solid lines in the deformed configuration. To first-order accuracy in components of $[\partial u / \partial X]$, the extensional strains of these fibres are $\epsilon_{11} = \partial u_1 / \partial X_1$ and $\epsilon_{22} = \partial u_2 / \partial X_2$, and the reduction of the angle between them is $\gamma_{12} = \partial u_2 / \partial X_1 + \partial u_1 / \partial X_2$. For the shear strain denoted ϵ_{12} , however, half of γ_{12} is used. Thus, considering all extensional and shear strains associated with infinitesimal fibres in the 1, 2, and 3 directions at a point of the material, the set of strains is given by

$$\epsilon_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial X_j} + \frac{\partial u_j}{\partial X_i} \right) \quad (i, j = 1, 2, 3). \quad (114)$$

The ϵ_{ij} are symmetric—*i.e.*, $\epsilon_{ij} = \epsilon_{ji}$ —and form a second-rank tensor (that is, if Cartesian reference axes 1', 2', and 3' were chosen instead and the ϵ_{kl}' were determined, then the ϵ_{kl}' are related to the ϵ_{ij} by the same equations that relate the stresses σ_{kl}' to the σ_{ij}). These mathematical features require that there exist principal strain directions; at every point of the continuum it is possible to identify three mutually perpendicular directions along which there is purely extensional strain, with no shear strain between these special directions. The directions are the principal strain directions, and the corresponding strains include the least and greatest extensional strains experienced by fibres through the material point considered. Invariants of the

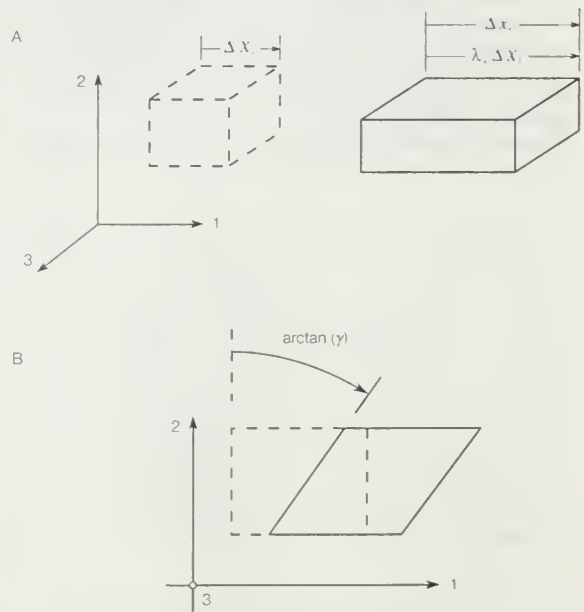


Figure 41: (A) Extensional strain and (B) simple shear strain, where the element drawn with dashed lines represents the reference configuration and the element drawn with solid lines represents the deformed configuration.

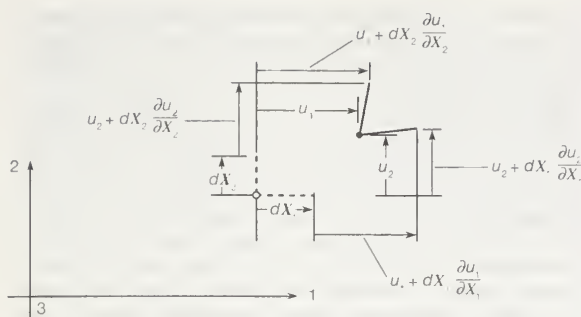


Figure 42: Relations of strains to gradients of displacement (see text).

strain tensor may be defined in a way paralleling those for the stress tensor.

An important fact to note is that the strains cannot vary in an arbitrary manner from point to point in the body. This is because the six strain components are all derivable from three displacement components. Restrictions on strain resulting from such considerations are called compatibility relations; the body would not fit together after deformation unless they were satisfied. Consider, for example, a state of plane strain in the 1, 2 plane (so that $\epsilon_{33} = \epsilon_{23} = \epsilon_{31} = 0$). The nonzero strains ϵ_{11} , ϵ_{22} , and ϵ_{12} cannot vary arbitrarily from point to point but must satisfy $\partial^2 \epsilon_{22} / \partial X_1^2 + \partial^2 \epsilon_{11} / \partial X_2^2 = 2\partial^2 \epsilon_{12} / \partial X_1 \partial X_2$, as may be verified by directly inserting the relations for strains in terms of displacements.

When the smallness of stretch and rotation of line elements allows use of the infinitesimal strain tensor, a derivative $\partial/\partial X_i$ will be very nearly identical to $\partial/\partial x_i$. Frequently, but not always, it will then be acceptable to ignore the distinction between the deformed and undeformed configurations in writing the governing equations of solid mechanics. For example, the differential equations of motion in terms of stress are rigorously correct only with derivatives relative to the deformed configuration, but, in the circumstances considered, the equations of motion can be written relative to the undeformed configuration. This is what is done in the most widely used variant of solid mechanics, in the form of the theory of linear elasticity. The procedure can be unsatisfactory and go badly wrong in some important cases, however, such as for columns that buckle under compressive loadings or for elastic-plastic materials when the slope of the stress versus strain relation is of the same order as existing stresses. Cases such as these are instead best approached through finite deformation theory.

Finite deformation and strain tensors. In the theory of finite deformations, extension and rotations of line elements are unrestricted as to size. For an infinitesimal fibre that deforms from an initial point given by the vector dX to the vector dx in the time t , the deformation gradient is defined by $F_{ij} = \partial x_i(X, t) / \partial X_j$; the 3×3 matrix $[F]$, with components F_{ij} , may be represented as a pure deformation, characterized by a symmetric matrix $[U]$, followed by a rigid rotation $[R]$. This result is called the polar decomposition theorem and takes the form, in matrix notation, $[F] = [R][U]$. For an arbitrary deformation, there exist three mutually orthogonal principal stretch directions at each point of the material; call these directions in the reference configuration $N^{(1)}$, $N^{(2)}$, $N^{(3)}$, and let the stretch ratios be λ_1 , λ_2 , λ_3 . Fibres in these three principal strain directions undergo extensional strain but have no shearing between them. Those three fibres in the deformed configuration remain orthogonal but are rotated by the operation $[R]$.

As noted earlier, an extensional strain may be defined by $E = g(\lambda)$, where $g(1) = 0$ and $g'(1) = 1$, with examples for $g(\lambda)$ given above. A finite strain tensor E_{ij} may then be defined based on any particular function $g(\lambda)$ by $E_{ij} = g(\lambda_1)N_i^{(1)}N_j^{(1)} + g(\lambda_2)N_i^{(2)}N_j^{(2)} + g(\lambda_3)N_i^{(3)}N_j^{(3)}$. Usually, it is rather difficult to actually solve for the λ 's and N 's associated with any general $[F]$, so it is not easy to use this strain definition. However, for the special choice identified as $g^H(\lambda) = (\lambda^2 - 1)/2$ above, it may be shown that $2E_{ij}^H =$

$$\sum_{k=1}^3 F_{ki} F_{kj} - \delta_{ij} = \partial u_i / \partial X_j + \partial u_j / \partial X_i + \sum_{k=1}^3 (\partial u_k / \partial X_i) (\partial u_k / \partial X_j),$$

which, like the finite strain generated by any other $g(\lambda)$, reduces to ϵ_{ij} when linearized in $[\partial u/\partial X]$.

STRESS-STRAIN RELATIONS

Linear elastic isotropic solid. The simplest type of stress-strain relation is that of the linear elastic solid, considered in circumstances for which $|\partial u_i / \partial X_j| \ll 1$ and for isotropic materials, whose mechanical response is independent of the direction of stressing. If a material point sustains a stress state $\sigma_{11} = \sigma$, with all other $\sigma_{ij} = 0$, it is subjected to uniaxial tensile stress. This can be realized in a homogeneous bar loaded by an axial force. The resulting strain may be rewritten as $\epsilon_{11} = \sigma/E$, $\epsilon_{22} = \epsilon_{33} = -\nu \epsilon_{11} = -\nu \sigma/E$, $\epsilon_{12} = \epsilon_{23} = \epsilon_{31} = 0$. Two new parameters have been introduced here, E and ν . E is called Young's modulus, and it has dimensions of [force]/[length]² and is measured in units such as the pascal (1 Pa = 1 N/m²), dyne/cm², or pounds per square inch (psi); ν , which equals the ratio of lateral strain to axial strain, is dimensionless and is called the Poisson ratio.

If the isotropic solid is subjected only to shear stress τ —i.e., $\sigma_{12} = \sigma_{21} = \tau$, with all other $\sigma_{ij} = 0$ —then the response is shearing strain of the same type, $\epsilon_{12} = \tau/2G$, $\epsilon_{23} = \epsilon_{31} = \epsilon_{11} = \epsilon_{22} = \epsilon_{33} = 0$. Notice that because $2\epsilon_{12} = \gamma_{12}$, this is equivalent to $\gamma_{12} = \tau/G$. The constant G introduced is called the shear modulus. (Frequently, the symbol μ is used instead of G .) The shear modulus G is not independent of E and ν but is related to them by $G = E/2(1 + \nu)$, as follows from the tensor nature of stress and strain. The general stress-strain relations are then

$$\epsilon_{ij} = (1 + \nu) \frac{\sigma_{ij}}{E} - \nu \delta_{ij} \frac{\sigma_{11} + \sigma_{22} + \sigma_{33}}{E} \quad (i, j = 1, 2, 3), \quad (115)$$

where δ_{ij} is defined as 1 when its indices agree and 0 otherwise.

These relations can be inverted to read $\sigma_{ij} = \lambda \delta_{ij} (\epsilon_{11} + \epsilon_{22} + \epsilon_{33}) + 2\mu \epsilon_{ij}$, where μ has been used rather than G as the notation for the shear modulus, following convention, and where $\lambda = 2\nu\mu/(1 - 2\nu)$. The elastic constants λ and μ are sometimes called the Lamé constants. Since ν is typically in the range $1/4$ to $1/3$ for hard polycrystalline solids, λ falls often in the range between μ and 2μ . (Navier's particle model with central forces leads to $\lambda = \mu$ for an isotropic solid.)

Another elastic modulus often cited is the bulk modulus K , defined for a linear solid under pressure p ($\sigma_{11} = \sigma_{22} = \sigma_{33} = -p$) such that the fractional decrease in volume is p/K . For example, consider a small cube of side length L in the reference state. If the length along, say, the 1 direction changes to $(1 + \epsilon_{11})L$, the fractional change of volume is $(1 + \epsilon_{11})(1 + \epsilon_{22})(1 + \epsilon_{33}) - 1 = \epsilon_{11} + \epsilon_{22} + \epsilon_{33}$, neglecting quadratic and cubic order terms in the ϵ_{ij} compared to the linear, as is appropriate when using linear elasticity. Thus, $K = E/3(1 - 2\nu) = \lambda + 2\mu/3$.

Thermal strains. Temperature change can also cause strain. In an isotropic material the thermally induced extensional strains are equal in all directions, and there are no shear strains. In the simplest cases, these thermal strains can be treated as being linear in the temperature change $\theta - \theta_0$ (where θ_0 is the temperature of the reference state), writing $\epsilon_{ij}^{\text{thermal}} = \delta_{ij} \alpha (\theta - \theta_0)$ for the strain produced by temperature change in the absence of stress. Here α is called the coefficient of thermal expansion. Thus, in cases of temperature change, ϵ_{ij} is replaced in the stress-strain relations above with $\epsilon_{ij} - \epsilon_{ij}^{\text{thermal}}$, with the thermal part given as a function of temperature. Typically, when temperature changes are modest, the small dependence of E and ν on temperature can be neglected.

Anisotropy. Anisotropic solids also are common in nature and technology. Examples are single crystals; polycrystals in which the grains are not completely random in their crystallographic orientation but have a "texture," typically owing to some plastic or creep flow process that has left a preferred grain orientation; fibrous biological materials such as wood or bone; and composite materials that, on a

microscale, either have the structure of reinforcing fibres in a matrix, with fibres oriented in a single direction or in multiple directions (e.g., to ensure strength along more than a single direction), or have the structure of a lamination of thin layers of separate materials. In the most general case, the application of any of the six components of stress induces all six components of strain, and there is no shortage of elastic constants. There would seem to be $6 \times 6 = 36$ in the most general case, but, as a consequence of the laws of thermodynamics, the maximum number of independent elastic constants is 21 (compared with 2 for isotropic solids). In many cases of practical interest, symmetry considerations reduce the number to far below 21. For example, crystals of cubic symmetry, such as rock salt (NaCl); face-centred cubic metals, such as aluminum, copper, or gold; body-centred cubic metals, such as iron at low temperatures or tungsten; and such nonmetals as diamond, germanium, or silicon have only three independent elastic constants. Solids with a special direction, and with identical properties along any direction perpendicular to that direction, are called transversely isotropic; they have five independent elastic constants. Examples are provided by fibre-reinforced composite materials, with fibres that are randomly emplaced but aligned in a single direction in an isotropic or transversely isotropic matrix, and by single crystals of hexagonal close packing such as zinc.

Transversely isotropic solids

General linear elastic stress-strain relations have the form

$$\sigma_{ij} = \sum_{k=1}^3 \sum_{l=1}^3 C_{ijkl} \epsilon_{kl}$$

where the coefficients C_{ijkl} are known as the tensor elastic moduli. Because the ϵ_{kl} are symmetric, one may choose $C_{ijkl} = C_{jilk}$, and, because the σ_{ij} are symmetric, $C_{ijkl} = C_{jikl}$. Hence the $3 \times 3 \times 3 \times 3 = 81$ components of C_{ijkl} reduce to the $6 \times 6 = 36$ mentioned. In cases of temperature change, the ϵ_{ij} above is replaced by $\epsilon_{ij} - \epsilon_{ij}^{thermal}$, where $\epsilon_{ij}^{thermal} = \alpha_{ij}(\theta - \theta_0)$ and α_{ij} is the set of thermal strain coefficients, with $\alpha_{ij} = \alpha_{ji}$. An alternative matrix notation is sometimes employed, especially in the literature on single crystals. That approach introduces 6-element columns of stress and strain $\{\sigma\}$ and $\{\epsilon\}$, defined so that the columns, when transposed (superscript T) or laid out as rows, are $\{\sigma\}^T = (\sigma_{11}, \sigma_{22}, \sigma_{33}, \sigma_{12}, \sigma_{23}, \sigma_{31})$ and $\{\epsilon\}^T = (\epsilon_{11}, \epsilon_{22}, \epsilon_{33}, 2\epsilon_{12}, 2\epsilon_{23}, 2\epsilon_{31})$. These forms assure that the scalar $\{\sigma\}^T \{d\epsilon\} = \text{tr}\{\sigma\} \{d\epsilon\}$ is an increment of stress working per unit volume. The stress-strain relations are then written $\{\sigma\} = [c]\{\epsilon\}$, where $[c]$ is the 6×6 matrix of elastic moduli. Thus, $c_{13} = C_{1133}$, $c_{15} = C_{1123}$, $c_{44} = C_{1212}$, and so on.

Thermodynamic considerations. In thermodynamic terminology, a state of purely elastic material response corresponds to an equilibrium state, and a process during which there is purely elastic response corresponds to a sequence of equilibrium states and hence to a reversible process. The second law of thermodynamics assures that the heat absorbed per unit mass can be written θds , where θ is the thermodynamic (absolute) temperature and s is the entropy per unit mass. Hence, writing the work per unit volume of reference configuration in a manner appropriate to cases when infinitesimal strain can be used, and letting ρ_0 be the density in that configuration, from the first law of thermodynamics it can be stated that $\rho_0 \theta ds + \text{tr}\{\sigma\} \{d\epsilon\} = \rho_0 de$, where e is the internal energy per unit mass. This relation shows that if e is expressed as a function of entropy s and strains $\{\epsilon\}$, and if e is written so as to depend identically on ϵ_{ij} and ϵ_{ji} , then $\sigma_{ij} = \rho_0 \partial e(\{\epsilon\}, s) / \partial \epsilon_{ij}$.

Helmholtz free energy

Alternatively, one may introduce the Helmholtz free energy f per unit mass, where $f = e - \theta s = f(\{\epsilon\}, \theta)$, and show that $\sigma_{ij} = \rho_0 \partial f(\{\epsilon\}, \theta) / \partial \epsilon_{ij}$. The latter form corresponds to the variables with which the stress-strain relations were written above. Sometimes $\rho_0 f$ is called the strain energy for states of isothermal (constant θ) elastic deformation; $\rho_0 e$ has the same interpretation for isentropic ($s = \text{constant}$) elastic deformation, achieved when the time scale is too short to allow heat transfer to or from a deforming element. Since the mixed partial derivatives must be independent of order, a consequence of the last equation is that $\partial \sigma_{ij}(\{\epsilon\}, \theta) / \partial \epsilon_{kl} = \partial \sigma_{kl}(\{\epsilon\}, \theta) / \partial \epsilon_{ij}$, which requires that $C_{ijkl} = C_{klij}$, or equivalently that the matrix $[c]$ be symmet-

ric, $[c] = [c]^T$, reducing the maximum possible number of independent elastic constraints from 36 to 21. The strain energy $W(\{\epsilon\})$ at constant temperature θ_0 is $W(\{\epsilon\}) \equiv \rho_0 f(\{\epsilon\}, \theta_0) = (1/2)\{\epsilon\}^T [c] \{\epsilon\}$.

The elastic moduli for isentropic response are slightly different from those for isothermal response. In the case of the isotropic material, it is convenient to give results in terms of G and K , the isothermal shear and bulk moduli. The isentropic moduli \bar{G} and \bar{K} are then $\bar{G} = G$ and $\bar{K} = K(1 + 9\theta_0 K \alpha^2 / \rho_0 c_e)$, where $c_e = \theta_0 \partial s(\{\epsilon\}, \theta) / \partial \theta$, evaluated at $\theta = \theta_0$ and $\{\epsilon\} = [0]$, is the specific heat at constant strain. The fractional change in the bulk modulus, given by the second term in the parentheses, is very small, typically on the order of 1 percent or less, even for metals and ceramics of relatively high α , on the order of 10^{-5} /kelvin.

The fractional change in absolute temperature during an isentropic deformation is found to involve the same small parameter: $[(\theta - \theta_0) / \theta_0]_{s = \text{const}} = -(9\theta_0 K \alpha^2 / \rho_0 c_e) [(\epsilon_{11} + \epsilon_{22} + \epsilon_{33}) / 3\alpha\theta_0]$. Values of α for most solid elements and inorganic compounds are in the range of 10^{-6} to 4×10^{-5} /kelvin; room temperature is about 300 kelvins, so $3\alpha\theta_0$ is typically in the range 10^{-3} to 4×10^{-2} . Thus, if the fractional change in volume is on the order of 1 percent, which is quite large for a metal or ceramic deforming in its elastic range, the fractional change in absolute temperature is also on the order of 1 percent. For those reasons, it is usually appropriate to neglect the alteration of the temperature field due to elastic deformation and hence to use purely mechanical formulations of elasticity in which distinctions between isentropic and isothermal response are neglected.

Finite elastic deformations. When elastic response under arbitrary deformation gradients is considered—because rotations, if not strains, are large or, in a material such as rubber, because the strains are large too—it is necessary to dispense with the infinitesimal strain theory. In such cases, the combined first and second laws of thermodynamics have the form $\rho_0 \theta ds + \det[F] \text{tr}\{[F]^{-1}[\sigma][dF]\} = \rho_0 de$, where $[F]^{-1}$ is the matrix inverse of the deformation gradient $[F]$. If a parcel of material is deformed by $[F]$ and then given some additional rigid rotation, the free energy f must be unchanged in that rotation. In terms of the polar decomposition $[F] = [R][U]$, this is equivalent to saying that f is independent of the rotation part $[R]$ of $[F]$, which is then equivalent to saying that f is a function of the finite strain measure $[E^M] = (1/2)\{[F]^T[F] - [I]\}$ based on change of metric or, for that matter, on any member of the family of material strain tensors. Thus,

$$\sigma_{ij} = (1/\det[F]) \sum_{k=1}^3 \sum_{l=1}^3 F_{ik} F_{jl} S_{kl}([E^M], \theta)$$

where $S_{kl} (= S_{lk})$ is sometimes called the second Piola-Kirchhoff stress and is given by $S_{kl} = \rho_0 \partial f([E^M], \theta) / \partial E_{kl}^M$, it being assumed that f has been written so as to have identical dependence on E_{kl}^M and E_{lk}^M .

Piola-Kirchhoff stress

Inelastic response. The above mode of expressing $[\sigma]$ in terms of $[S]$ is valid for solids showing viscoelastic or plastic response as well, except that $[S]$ is then to be regarded not only as a function of the present $[E^M]$ and θ but also as dependent on the prior history of both. Assuming that such materials show elastic response to sudden stress changes or to small unloading from a plastically deforming state, $[S]$ may still be expressed as a derivative of f , as above, but the derivative is understood as being taken with respect to an elastic variation of strain and is to be taken at fixed θ and with fixed prior inelastic deformation and temperature history. Such dependence on history is sometimes represented as a dependence of f on internal state variables whose laws of evolution are part of the inelastic constitutive description. There are also simpler models of inelastic response, and the most commonly employed forms for plasticity and creep in isotropic solids are presented next.

To a good approximation, plastic deformation of crystalline solids causes no change in volume; and hydrostatic changes in stress, amounting to equal change of all normal stresses, have no effect on plastic flow, at least for changes that are of the same order or magnitude as the strength

of the solid in shear. Thus, plastic response is usually formulated in terms of deviatoric stress, which is defined by $\tau_{ij} = \sigma_{ij} - \delta_{ij}(\sigma_{11} + \sigma_{22} + \sigma_{33})/3$. Following Richard von Mises, in a procedure that is found to agree moderately well with experiment, the plastic flow relation is formulated in terms of the second invariant of deviatoric stress, commonly rewritten as $\bar{\sigma} = \sqrt{(3/2)\text{tr}([\tau][\tau])}$ and called the

Equivalent tensile stress

equivalent tensile stress. The definition is made so that, for a state of uniaxial tension, $\bar{\sigma}$ equals the tensile stress, and the stress-strain relation for general stress states is formulated in terms of data from the tensile test. In particular, a plastic strain $\bar{\epsilon}^p$ in a uniaxial tension test is defined from $\bar{\epsilon}^p = \bar{\epsilon} - \bar{\sigma}/E$, where $\bar{\epsilon}$ is interpreted as the strain in the tensile test according to the logarithmic definition $\bar{\epsilon} = \ln \lambda$, the elastic modulus E is assumed to remain unchanged with deformation, and $\bar{\sigma}/E \ll 1$.

Thus, in the rate-independent plasticity version of the theory, tensile data (or compressive, with appropriate sign reversals) from a monotonic load test is assumed to define a function $\bar{\epsilon}^p(\bar{\sigma})$. In the viscoplastic or high-temperature creep versions of the theory, tensile data is interpreted to define $d\bar{\epsilon}^p/dt$ as a function of $\bar{\sigma}$ in the simplest case, representing, for example, secondary creep, and as a function of $\bar{\sigma}$ and $\bar{\epsilon}^p$ in theories intended to represent transient creep effects or rate-sensitive response at lower temperatures. Consider first the rigid-plastic material model in which elastic deformability is ignored altogether, as is sometimes appropriate for problems of large plastic flow, as in metal forming or long-term creep in the Earth's mantle or for analysis of plastic collapse loads on structures. The rate of deformation tensor D_{ij} is defined by $2D_{ij} = \partial v_i/\partial x_j + \partial v_j/\partial x_i$, and in the rigid-plastic case $[D]$ can be equated to what may be considered its plastic part $[D^p]$, given as $D_{ij}^p = 3(d\bar{\epsilon}^p/dt)\tau_{ij}/2\bar{\sigma}$. The numerical factors secure agreement between D_{ij}^p and $d\bar{\epsilon}^p/dt$ for uniaxial tension in the 1-direction. Also, the equation implies that $D_{11}^p + D_{22}^p + D_{33}^p = 0$ and that $d\bar{\epsilon}^p/dt = \sqrt{(2/3)\text{tr}([D^p][D^p])}$, which must be

integrated over previous history to get $\bar{\epsilon}^p$ as required for viscoplastic models in which $d\bar{\epsilon}^p/dt$ is a function of $\bar{\sigma}$ and $\bar{\epsilon}^p$. In the rate-independent version, $[D^p]$ is defined as zero whenever $\bar{\sigma}$ is less than the highest value that it has attained in the previous history or when the current value of $\bar{\sigma}$ is the highest value but $d\bar{\sigma}/dt < 0$. (In the elastic-plastic context, this means that "unloading" involves only elastic response.) For the ideally plastic solid, which is idealized to be able to flow without increase of stress when $\bar{\sigma}$ equals the yield strength level, $d\bar{\epsilon}^p/dt$ is regarded as an undetermined but necessarily nonnegative parameter, which can be determined (sometimes not uniquely) only through the complete solution of a solid mechanics boundary-value problem.

The elastic-plastic material model is then formulated by writing $D_{ij} = D_{ij}^e + D_{ij}^p$, where D_{ij}^e is given in terms of stress and possibly stress rate as above and where the elastic deformation rates $[D^e]$ are related to stresses by the usual linear elastic expression $D_{ij}^e = (1 + \nu)\dot{\sigma}_{ij}/E - \nu\delta_{ij}(\dot{\sigma}_{11} + \dot{\sigma}_{22} + \dot{\sigma}_{33})/E$. Here the stress rates are expressed as the Jaumann co-rotational rates $\dot{\sigma}_{ij}^* = \dot{\sigma}_{ij} + \sum_{k=1}^3 (\sigma_{ik}\Omega_{kj} - \Omega_{ik}\sigma_{kj})$, where $\dot{\sigma}_{ij} = d\sigma_{ij}/dt$ is a derivative

following the motion of a material point and where the spin Ω_{ij} is defined by $2\Omega_{ij} = \partial v_j/\partial x_i - \partial v_i/\partial x_j$. The co-rotational stress rates are those calculated by an observer who spins with the average angular velocity of a material element. The elastic part of the stress-strain relation should be consistent with the existence of a free energy f , as discussed above. This is not strictly satisfied by the form just given, but the differences between it and one which is consistent in that way involves additional terms that are on the order of $\bar{\sigma}/E^2$ times the $\dot{\sigma}_{kl}^*$ and are negligible in typical cases in which the theory is used, since $\bar{\sigma}/E$ is usually an extremely small fraction of unity, say, 10^{-4} to 10^{-2} . A small-strain version of the theory is in common use for purposes of elastic-plastic stress analysis. In these cases, $[D]$ is replaced by $\partial[e(X, t)]/\partial t$, where $[e]$ is the small-strain tensor, $\partial/\partial x$ with $\partial/\partial X$ in all equations, and

Small-strain theory

$[\dot{\sigma}^*]$ with $\partial[\sigma(X, t)]/\partial t$. The last two steps cannot always be justified, even in cases of very small strain when, for example, in a rate-independent material, $d\bar{\sigma}/d\bar{\epsilon}^p$ is not large compared to $\bar{\sigma}$ or when rates of rotation of material fibres can become much larger than rates of stretching, which is a concern for buckling problems even in purely elastic solids.

PROBLEMS INVOLVING ELASTIC RESPONSE

Equations of motion of linear elastic bodies. The final equations of the purely mechanical theory of linear elasticity (*i.e.*, when coupling with the temperature field is neglected, or when either isothermal or isentropic response is assumed) are obtained as follows. The stress-strain relations are used, and the strains are written in terms of displacement gradients. The final expressions for stress are inserted into the equations of motion, replacing $\partial/\partial x$ with $\partial/\partial X$ in those equations. In the case of an isotropic and homogenous solid, these reduce to

$$(\lambda + \mu)\nabla(\nabla \cdot \mathbf{u}) + \mu\nabla^2\mathbf{u} + \mathbf{f} = \rho\frac{\partial^2\mathbf{u}}{\partial t^2}, \quad (116)$$

known as the Navier equations (here, $\nabla = \mathbf{e}_1\partial/\partial X_1 + \mathbf{e}_2\partial/\partial X_2 + \mathbf{e}_3\partial/\partial X_3$, and ∇^2 is the Laplacian operator defined by $\nabla \cdot \nabla$, or $\partial^2/\partial x_1^2 + \partial^2/\partial x_2^2 + \partial^2/\partial x_3^2$, and, as described earlier, λ and μ are the Lamé constants, \mathbf{u} the displacement, \mathbf{f} the body force, and ρ the density of the material). Such equations hold in the region V occupied by the solid; on the surface S one prescribes each component of \mathbf{u} , or each component of the stress vector \mathbf{T} (expressed in terms of $[\partial u/\partial X]$), or sometimes mixtures of components or relations between them. For example, along a freely slipping planar interface with a rigid solid, the normal component of \mathbf{u} and the two tangential components of \mathbf{T} would be prescribed, all as zero.

Navier equations

Body wave solutions. By looking for body wave solutions in the form $\mathbf{u}(X, t) = \mathbf{p}f(\mathbf{n} \cdot \mathbf{X} - ct)$, where unit vector \mathbf{n} is the propagation direction, \mathbf{p} is the polarization, or direction of particle motion, and c is the wave speed, one may show for the isotropic material that solutions exist for arbitrary functions f if either $c = c_d \equiv \sqrt{(\lambda + 2\mu)/\rho}$ and $\mathbf{p} = \mathbf{n}$, or $c = c_s \equiv \sqrt{\mu/\rho}$ and $\mathbf{p} \cdot \mathbf{n} = 0$. The first case, with particle displacements in the propagation direction, describes longitudinal, or dilatational, waves; and the latter case, which corresponds to two linearly independent displacement directions, both transverse to the propagation direction, describes transverse, or shear, waves.

Linear elastic beam. The case of a beam treated as a linear elastic line may also be considered. Let the line along the 1-axis (see Figure 43), have properties that are uniform along its length and have sufficient symmetry that bending it by applying a torque about the 3-direction causes the line to deform into an arc lying in the 1,2-plane. Make an imaginary cut through the line, and let the forces and torque acting at that section on the part lying in the direction of decreasing X_1 be denoted as a shear force V in the positive 2-direction, an axial force P in the positive 1-direction, and torque M , commonly called a bending moment, about the positive 3-direction. The linear and angular momentum principles then require that the actions at that section on the part of the line lying along the direction of increasing X_1 be of equal magnitude but opposite sign.

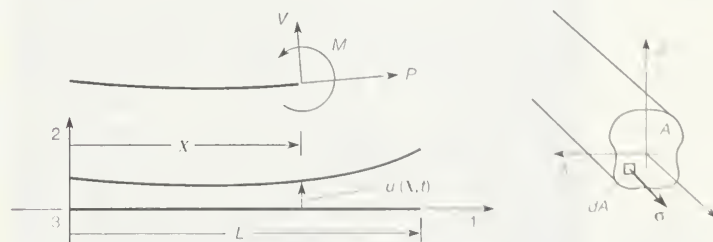


Figure 43: Transverse motion of an initially straight beam, shown at left as an elastic line and at right as a solid of finite section (see text).

Now let the line be loaded by transverse force F per unit length, directed in the 2-direction, and make assumptions on the smallness of deformation consistent with those of linear elasticity. Let ρA be the mass per unit length (so that A can be interpreted as the cross-sectional area of a homogeneous beam of density ρ) and let u be the transverse displacement in the 2-direction. Then, writing X for X_1 , the linear and angular momentum principles require that $\partial V/\partial X + F = \rho A \partial^2 u/\partial t^2$ and $\partial M/\partial X + V = 0$, where rotary inertia has been neglected in the second equation, as is appropriate for disturbances which are of a wavelength that is long compared to cross-sectional dimensions. The curvature κ of the elastic line can be approximated by $\kappa = \partial^2 u/\partial X^2$ for the small deformation situation considered, and the equivalent of the stress-strain relation is to assume that κ is a function of M at each point along the line. The function can be derived by the analysis of stress and strain in pure bending and is $M = EI\kappa$, with the moment of inertia $I = \int_A (X_2)^2 dA$ for uniform elastic properties over all the cross section and with the 1-axis passing through the section centroid. Hence, the equation relating transverse load and displacement of a linear elastic beam is $-\partial^2(EI\partial^2 u/\partial X^2)/\partial X^2 + F = \rho A \partial^2 u/\partial t^2$, and this is to be solved subject to two boundary conditions at each end of the elastic line. Examples are $u = \partial u/\partial X = 0$ at a completely restrained ("built in") end, $u = M = 0$ at an end that is restrained against displacement but not rotation, and $V = M = 0$ at a completely unrestrained (free) end. The beam will be reconsidered later in an analysis of response with initial stress present.

Boundary conditions

The preceding derivation was presented in the spirit of the model of a beam as the elastic line of Euler. The same equations of motion may be obtained by the following five steps: (1) integrate the three-dimensional equations of motion over a section, writing $V = \int_A \sigma_{12} dA$; (2) integrate the product of X_2 and those equations over a section, writing $M = -\int_A X_2 \sigma_{11} dA$; (3) assume that planes initially perpendicular to fibres lying along the 1-axis remain perpendicular during deformation, so that $\epsilon_{11} = \epsilon_0(X, t) - X_2 \kappa(X, t)$, where $X \equiv X_1$, $\epsilon_0(X, t)$ is the strain of the fibre along the 1-axis, and $\kappa(X, t) = \partial^2 u/\partial X^2$, where $u(X, t)$ is u_2 for the fibre initially along the 1-axis; (4) assume that the stress σ_{11} relates to strain as if each point were under uniaxial tension, so that $\sigma_{11} = E\epsilon_{11}$; and (5) neglect terms of order h^2/L^2 compared to unity, where h is a typical cross-section dimension and L is a scale length for variations along the direction of the 1-axis. In step (1) the average of u_2 over area A enters but may be interpreted as the displacement u of step (3) to the order retained in (5). The kinematic assumption (3) together with (5), if implemented under conditions such that there are no loadings to generate a net axial force P , requires that $\epsilon_0(X, t) = 0$ and that $\kappa(X, t) = M(X, t)/EI$ when the 1-axis has been chosen to pass through the centroid of the cross section. Hence, according to these approximations, $\sigma_{11} = -X_2 M(X, t)/I = -X_2 E \partial^2 u(X, t)/\partial X^2$. The expression for σ_{11} is exact for static equilibrium under pure bending, since assumptions (3) and (4) are exact and (5) is then irrelevant. This motivates the use of assumptions (3) and (4) in a situation that does not correspond to pure bending.

Sometimes it is necessary to deal with solids that are already under stress in the reference configuration that is chosen for measuring strain. As a simple example, suppose that the beam just discussed is under an initial uniform tensile stress $\sigma_{11} = \sigma^0$ —that is, the axial force $P = \sigma^0 A$. If σ^0 is negative and of significant magnitude, one generally refers to the beam as a column; if it is large and positive, the beam might respond more like a taut string. The initial stress σ^0 contributes a term to the equations of small transverse motion, which now becomes $-\partial^2(EI\partial^2 u/\partial X^2)/\partial X^2 + \sigma^0 A \partial^2 u/\partial X^2 + F = \rho A \partial^2 u/\partial t^2$.

Free vibrations. Suppose that the beam is of length L , is of uniform properties, and is hinge-supported at its ends at $X = 0$ and $X = L$ so that $u = M = 0$ there. Then free transverse motions of the beam, solving the above equation with $F = 0$, are described by any linear combination of the real part of solutions that have the form $u = C_n \exp(i\omega_n t) \sin(n\pi X/L)$, where n is any positive integer, C_n is an arbitrary complex constant, and where

$$\rho A \omega_n^2 = \left(\frac{n\pi}{L}\right)^4 EI \left[1 + \left(\frac{\sigma^0}{E}\right) \left(\frac{AL^2}{n^2 \pi^2 I}\right)\right] \quad (117)$$

defines the angular vibration frequency ω_n , associated with the n th mode, in units of radians per unit time. The number of vibration cycles per unit time is $\omega_n/2\pi$. Equation (117) is arranged so that the term in the brackets shows the correction, from unity, of what would be the expression giving the frequencies of free vibration for a beam when there is no σ^0 . The correction from unity can be quite significant, even though σ^0/E is always much smaller than unity (for interesting cases, 10^{-6} to, say, 10^{-3} would be a representative range; few materials in bulk form would remain elastic or resist fracture at higher σ^0/E , although good piano wire could reach about 10^{-2}). The correction term's significance results because σ^0/E is multiplied by a term that can become enormous for a beam that is long compared to its thickness; for a square section of side length h , that term (at its largest, when $n = 1$) is $AL^2/\pi^2 I \approx 1.2L^2/h^2$, which can combine with a small σ^0/E to produce a correction term within the brackets that is quite non-negligible compared to unity. When $\sigma^0 > 0$ and L is large enough to make the bracketed expression much larger than unity, the EI term cancels out and the beam simply responds like a stretched string (here, string denotes an object that is unable to support a bending moment). When the vibration mode number n is large enough, however, the stringlike effects become negligible and beamlike response takes over; at sufficiently high n that L/n is reduced to the same order as h , the simple beam theory becomes inaccurate and should be replaced by three-dimensional elasticity or, at least, an improved beam theory that takes into account rotary inertia and shear deformability. (While the option of using three-dimensional elasticity for such a problem posed an insurmountable obstacle over most of the history of the subject, by 1990 the availability of computing power and easily used software reduced it to a routine problem that could be studied by an undergraduate engineer or physicist using the finite-element method or some other computational mechanics technique.)

Buckling. An important case of compressive loading is that in which $\sigma^0 < 0$, which can lead to buckling. Indeed, if $\sigma^0 A < -\pi^2 EI/L^2$, then the ω_n^2 is negative, at least for $n = 1$, which means that the corresponding ω_n is of the form $\pm ib$, where b is a positive real number, so that the $\exp(i\omega_n t)$ term has a time dependence of a type that no longer involves oscillation but, rather, exponential growth, $\exp(bt)$. The critical compressive force, $\pi^2 EI/L^2$, that causes this type of behaviour is called the Euler buckling load; different numerical factors are obtained for different end conditions. The acceleration associated with the $n = 1$ mode becomes small in the vicinity of the critical load and vanishes at that load. Thus solutions are possible, at the buckling load, for which the column takes a deformed shape without acceleration; for that reason, an approach to buckling problems that is equivalent for what, in dynamic terminology, are called conservative systems is to seek the first load at which an alternate equilibrium solution $u = u(X)$, other than $u = 0$, may exist.

Instability by divergence—that is, with growth of displacement in the form $\exp(bt)$ —is representative of conservative systems. Columns under nonconservative loadings by, for example, a follower force, which has the property that its line of action rotates so as to be always tangent to the beam centerline at its place of application, can exhibit a flutter instability in which the dynamic response is proportional to the real or imaginary part of a term such as $\exp(iat)\exp(bt)$ —i.e., an oscillation with exponentially growing amplitude. Such instabilities also arise in the coupling between fluid flow and elastic structural response, as in the subfield called aeroelasticity. The prototype is the flutter of an airplane wing—that is, a torsional oscillation of the wing, of growing amplitude, which is driven by the coupling between rotation of the wing and the development of aerodynamic forces related to the angle of attack; the coupling feeds more energy into the structure with each cycle.

Of course, instability models that are based on linearized theories and predicting exponential growth in time actu-

Instability

Initial stress

ally reveal no more than that the system is deforming out of the range for which the mathematical model applies. Proper nonlinear theories that take account of the finiteness of rotation, and sometimes the large and possibly nonelastic strain of material fibres, are necessary to really understand the phenomena. An important subclass of such nonlinear analyses for conservative systems involves the static post-buckling response of a perfect structure, such as a perfectly straight column or perfectly spherical shell. That post-buckling analysis allows one to determine

if increasing force is required for very large displacement to develop during the buckle or whether the buckling is of a more highly unstable type for which the load must diminish with buckling amplitude in order to still satisfy the equilibrium equations. The latter type of behaviour describes a structure whose maximum load (that is, the largest load it can support without collapsing) shows strong sensitivity to very small imperfections of material or geometry, as is the case with many shell structures.

(J.R.R.)

FLUID MECHANICS

Fluid mechanics is concerned with the response of fluids to forces exerted upon them. It is a branch of classical physics with applications of great importance in hydraulic and aeronautical engineering, chemical engineering, meteorology, and zoology.

The most familiar fluid is of course water, and an encyclopaedia of the 19th century probably would have dealt with the subject under the separate headings of hydrostatics, the science of water at rest, and hydrodynamics, the science of water in motion. Archimedes founded hydrostatics in about 250 BC when, according to legend, he leapt out of his bath and ran naked through the streets of Syracuse crying "Eureka!"; it has undergone rather little development since. The foundations of hydrodynamics, on the other hand, were not laid until the 18th century when mathematicians such as Leonhard Euler and Daniel Bernoulli began to explore the consequences, for a virtually continuous medium like water, of the dynamic principles that Newton had enunciated for systems composed of discrete particles. Their work was continued in the 19th century by several mathematicians and physicists of the first rank, notably G.G. Stokes and William Thomson. By the end of the century explanations had been found for a host of intriguing phenomena having to do with the flow of water through tubes and orifices, the waves that ships moving through water leave behind them, raindrops on windowpanes, and the like. There was still no proper understanding, however, of problems as fundamental as that of water flowing past a fixed obstacle and exerting a drag force upon it; the theory of potential flow, which worked so well in other contexts, yielded results that at relatively high flow rates were grossly at variance with experiment. This problem was not properly understood until 1904, when the German physicist Ludwig Prandtl introduced the concept of the boundary layer (see below *Hydrodynamics: Boundary layers and separation*). Prandtl's career continued into the period in which the first manned aircraft were developed. Since that time, the flow of air has been of as much interest to physicists and engineers as the flow of water, and hydrodynamics has, as a consequence, become fluid dynamics. The term fluid mechanics, as used here, embraces both fluid dynamics and the subject still generally referred to as hydrostatics.

One other representative of the 20th century who deserves mention here besides Prandtl is Geoffrey Taylor of England. Taylor remained a classical physicist while most of his contemporaries were turning their attention to the problems of atomic structure and quantum mechanics, and he made several unexpected and important discoveries in the field of fluid mechanics. The richness of fluid mechanics is due in large part to a term in the basic equation of the motion of fluids which is nonlinear—*i.e.*, one that involves the fluid velocity twice over. It is characteristic of systems described by nonlinear equations that under certain conditions they become unstable and begin behaving in ways that seem at first sight to be totally chaotic. In the case of fluids, chaotic behaviour is very common and is called turbulence. Mathematicians have now begun to recognize patterns in chaos that can be analyzed fruitfully, and this development suggests that fluid mechanics will remain a field of active research well into the 21st century. (For a discussion of the concept of chaos, see *PHYSICAL SCIENCE, PRINCIPLES OF.*)

Fluid mechanics is a subject with almost endless ramifi-

cations, and the account that follows is necessarily incomplete. Some knowledge of the basic properties of fluids will be needed; a survey of the most relevant properties is given in the next section. For further details, see *THERMODYNAMICS, PRINCIPLES OF*; and *MATTER: Liquid state*.

Basic properties of fluids

Fluids are not strictly continuous media in the way that all the successors of Euler and Bernoulli have assumed, for they are composed of discrete molecules. The molecules, however, are so small and, except in gases at very low pressures, the number of molecules per millilitre is so enormous that they need not be viewed as individual entities. There are a few liquids, known as liquid crystals, in which the molecules are packed together in such a way as to make the properties of the medium locally anisotropic, but the vast majority of fluids (including air and water) are isotropic. In fluid mechanics, the state of an isotropic fluid may be completely described by defining its mean mass per unit volume, or density (ρ), its temperature (T), and its velocity (\mathbf{v}) at every point in space, and just what the connection is between these macroscopic properties and the positions and velocities of individual molecules is of no direct relevance.

A word perhaps is needed about the difference between gases and liquids, though the difference is easier to perceive than to describe. In gases the molecules are sufficiently far apart to move almost independently of one another, and gases tend to expand to fill any volume available to them. In liquids the molecules are more or less in contact, and the short-range attractive forces between them make them cohere; the molecules are moving too fast to settle down into the ordered arrays that are characteristic of solids, but not so fast that they can fly apart. Thus, samples of liquid can exist as drops or as jets with free surfaces, or they can sit in beakers constrained only by gravity, in a way that samples of gas cannot. Such samples may evaporate in time, as molecules one by one pick up enough speed to escape across the free surface and are not replaced. The lifetime of liquid drops and jets, however, is normally long enough for evaporation to be ignored.

There are two sorts of stress that may exist in any solid or fluid medium, and the difference between them may be illustrated by reference to a brick held between two hands. If the holder moves his hands toward each other, he exerts pressure on the brick; if he moves one hand toward his body and the other away from it, then he exerts what is called a shear stress. A solid substance such as a brick can withstand stresses of both types, but fluids, by definition, yield to shear stresses no matter how small these stresses may be. They do so at a rate determined by the fluid's viscosity. This property, about which more will be said later, is a measure of the friction that arises when adjacent layers of fluid slip over one another. It follows that the shear stresses are everywhere zero in a fluid at rest and in equilibrium, and from this it follows that the pressure (that is, force per unit area) acting perpendicular to all planes in the fluid is the same irrespective of their orientation (Pascal's law). For an isotropic fluid in equilibrium there is only one value of the local pressure (p) consistent with the stated values for ρ and T . These three quantities are linked together by what is called the equation of state for the fluid.

Equation of state

For gases at low pressures the equation of state is simple and well known. It is

$$p = \left(\frac{RT}{M} \right) \rho, \quad (118)$$

where R is the universal gas constant (8.3 joules per degree Celsius per mole) and M is the molar mass, or an average molar mass if the gas is a mixture; for air, the appropriate average is about 29×10^{-3} kilogram per mole. For other fluids knowledge of the equation of state is often incomplete. Except under very extreme conditions, however, all one needs to know is how the density changes when the pressure is changed by a small amount, and this is described by the compressibility of the fluid—either the isothermal compressibility, β_T , or the adiabatic compressibility, β_S , according to circumstance. When an element of fluid is compressed, the work done on it tends to heat it up. If the heat has time to drain away to the surroundings and the temperature of the fluid remains essentially unchanged throughout, then β_T is the relevant quantity. If virtually none of the heat escapes, as is more commonly the case in flow problems because the thermal conductivity of most fluids is poor, then the flow is said to be adiabatic, and β_S is needed instead. (The S refers to entropy, which remains constant in an adiabatic process provided that it takes place slowly enough to be treated as “reversible” in the thermodynamic sense.) For gases that obey equation (118), it is evident that p and ρ are proportional to one another in an isothermal process, and

$$\beta_T = \rho^{-1} \left(\frac{\partial \rho}{\partial p} \right)_T = \rho^{-1}. \quad (119)$$

In reversible adiabatic processes for such gases, however, the temperature rises on compression at a rate such that

$$T \propto \rho^{(\gamma-1)}, \quad p \propto \rho^\gamma, \quad (120)$$

and

$$\beta_S = \rho^{-1} \left(\frac{\partial \rho}{\partial p} \right)_S = (\gamma p)^{-1} = \frac{\beta_T}{\gamma}, \quad (121)$$

where γ is about 1.4 for air and takes similar values for other common gases. For liquids the ratio between the isothermal and adiabatic compressibilities is much closer to unity. For liquids, however, both compressibilities are normally much less than ρ^{-1} , and the simplifying assumption that they are zero is often justified.

The factor γ is not only the ratio between two compressibilities; it is also the ratio between two principal specific heats. The molar specific heat is the amount of heat required to raise the temperature of one mole through one degree. This is greater if the substance is allowed to expand as it is heated, and therefore to do work, than if its volume is fixed. The principal molar specific heats, C_p and C_v , refer to heating at constant pressure and constant volume, respectively, and

$$\gamma = \frac{C_p}{C_v}. \quad (122)$$

For air, C_p is about $3.5 R$.

Solids can be stretched without breaking, and liquids, though not gases, can withstand stretching, too. Thus, if the pressure is steadily reduced in a specimen of very pure water, bubbles will ultimately appear, but they may not do so until the pressure is negative and well below -10^7 newton per square metre; this is 100 times greater in magnitude than the (positive) pressure exerted by the Earth's atmosphere. Water owes its high ideal strength to the fact that rupture involves breaking links of attraction between molecules on either side of the plane on which rupture occurs; work must be done to break these links. However, its strength is drastically reduced by anything that provides a nucleus at which the process known as cavitation (formation of vapour- or gas-filled cavities) can begin, and a liquid containing suspended dust particles or dissolved gases is liable to cavitate quite easily.

Work also must be done if a free liquid drop of spherical shape is to be drawn out into a long thin cylinder or deformed in any other way that increases its surface area. Here again work is needed to break intermolecular links.

The surface of a liquid behaves, in fact, as if it were an elastic membrane under tension, except that the tension exerted by an elastic membrane increases when the membrane is stretched in a way that the tension exerted by a liquid surface does not. Surface tension is what causes liquids to rise up capillary tubes, what supports hanging liquid drops, what limits the formation of ripples on the surface of liquids, and so on.

Surface tension

Hydrostatics

It is common knowledge that the pressure of the atmosphere (about 10^5 newtons per square metre) is due to the weight of air above the Earth's surface, that this pressure falls as one climbs upward, and, correspondingly, that pressure increases as one dives deeper into a lake (or comparable body of water). Mathematically, the rate at which the pressure in a stationary fluid varies with height z in a vertical gravitational field of strength g is given by

$$\frac{dp}{dz} = -\rho g. \quad (123)$$

If ρ and g are both independent of z , as is more or less the case in lakes, then

$$p(z) = p(0) - \rho g z. \quad (124)$$

This means that, since ρ is about 10^3 kilograms per cubic metre for water and g is about 10 metres per second squared, the pressure is already twice the atmospheric value at a depth of 10 metres. Applied to the atmosphere, equation (124) would imply that the pressure falls to zero at a height of about 10 kilometres. In the atmosphere, however, the variation of ρ with z is far from negligible and (124) is unreliable as a consequence: a better approximation is given below in the section *Hydrodynamics: Compressible flow in gases*.

Instruments for comparing pressures are called differential manometers, and the simplest such instrument is a U-tube containing liquid, as shown in Figure 44A. The two pressures of interest, p_1 and p_2 , are transmitted to the two ends of the liquid column through an inert gas—the density of which is negligible by comparison with the liquid density, ρ —and the difference of height, h , of the two menisci is measured. It is a consequence of (124) that

$$p_1 - p_2 = \rho g h. \quad (125)$$

A barometer for measuring the pressure of the atmosphere in absolute terms is simply a manometer in which p_2 is made zero, or as close to zero as is feasible. The barometer invented in the 17th century by the Italian physicist and mathematician Evangelista Torricelli, and still in use today, is a U-tube that is sealed at one end (see Figure 44B). It may be filled with liquid, with the sealed end downward, and then inverted. On inversion, a negative pressure may momentarily develop at the top of the liquid column if the column is long enough; however, cavitation normally occurs there and the column falls away from the sealed end of the tube, as shown in the figure. Between the two exists what Torricelli thought of as a vacuum, though it may be very far from that condition if the barometer has been filled without scrupulous precautions to ensure that all dissolved or adsorbed gases, which would otherwise collect in this space, have first been removed. Even if no contaminating gas is present, the Torricellian vacuum always contains the vapour of the liquid, and this exerts a pressure which may be small but is never quite zero. The liquid conventionally used in a Torricelli barometer is of course mercury, which has a low vapour pressure and a high density. The high density means that h is only about 760 millimetres; if water were used, it would have to be about 10 metres instead.

Figure 44C illustrates the principle of the siphon. The top container is open to the atmosphere, and the pressure in it, p_2 , is therefore atmospheric. To balance this and the weight of the liquid column in between, the pressure p_1 in the bottom container ought to be greater by $\rho g h$. If the bottom container is also open to the atmosphere, then equilibrium is clearly impossible; the weight of the liquid column prevails and causes the liquid to flow downward.

Differential manometers

Molar specific heat

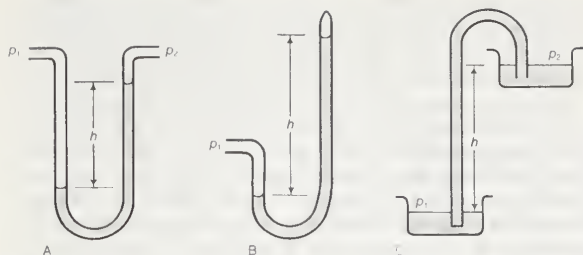


Figure 44: Schematic representations of (A) a differential manometer, (B) a Torricellian barometer, and (C) a siphon.

The siphon operates only as long as the column is continuous; it fails if a bubble of gas collects in the tube or if cavitation occurs. Cavitation therefore limits the level differences over which siphons can be used, and it also limits (to about 10 metres) the depth of wells from which water can be pumped using suction alone.

Consider now a cube of side d totally immersed in liquid with its top and bottom faces horizontal. The pressure on the bottom face will be higher than on the top by $\rho g d$, and, since pressure is force per unit area and the area of a cube face is d^2 , the resultant upthrust on the cube is $\rho g d^3$. This is a simple example of the so-called Archimedes' principle, which states that the upthrust experienced by a submerged or floating body is always equal to the weight of the liquid that the body displaces. As Archimedes must have realized, there is no need to prove this by detailed examination of the pressure difference between top and bottom. It is obviously true, whatever the body's shape. It is obvious because, if the solid body could somehow be removed and if the cavity thereby created could somehow be filled with more fluid instead, the whole system would then be in equilibrium. The extra fluid would, however, then be experiencing the upthrust previously experienced by the solid body, and it would not be in equilibrium unless this were just sufficient to balance its weight.

Archimedes' problem was to discover, by what would nowadays be called a nondestructive test, whether the crown of King Hieron II was made of pure gold or of gold diluted with silver. He understood that the pure metal and the alloy would differ in density and that he could determine the density of the crown by weighing it to find its mass and making a separate measurement of its volume. Perhaps the inspiration that struck him (in his bath) was that one can find the volume of any object by submerging it in liquid in something like a measuring cylinder (*i.e.*, in a container with vertical sides that have been suitably graduated) and measuring the displacement of the liquid surface. If so, he no doubt realized soon afterward that a more elegant and more accurate method for determining density can be based on the principle that bears his name. This method involves weighing the object twice, first, when it is suspended in a vacuum (suspension in air will normally suffice) and, second, when it is totally submerged in a liquid of density ρ . If the density of the object is ρ' , the ratio between the two weights must be

$$\frac{W_2}{W_1} = \frac{(\rho' - \rho)}{\rho'} \quad (126)$$

If ρ' is less than ρ , then W_2 , according to equation (126), is negative. What that means is that the object does not submerge of its own accord; it has to be pushed downward to make it do so. If an object with a mean density less than that of water is placed in a lake and not subjected to any downward force other than its own weight, it naturally floats on the surface, and Archimedes' principle shows that in equilibrium the volume of water which it displaces is a fraction ρ'/ρ of its own volume. A hydrometer is an object graduated in such a way that this fraction may be measured. By floating a hydrometer first in water of density ρ_0 and then in some other liquid of density ρ_1 and comparing the readings, one may determine the ratio ρ_1/ρ_0 —*i.e.*, the specific gravity of the other liquid.

In what orientation an object floats is a matter of grave concern to those who design boats and those who travel in them. A simple example will suffice to illustrate the factors that determine orientation. Figure 45 shows three

of the many possible orientations that a uniform square prism might adopt when floating, with half its volume submerged in a liquid for which $\rho = 2\rho'$; they are separated by rotations of 22.5° . In each of these diagrams, C is the centre of mass of the prism, and B, a point known as the centre of buoyancy, is the centre of mass of the displaced water. The distributed forces acting on the prism are equivalent to its weight acting downward through C and to the equal weight of the displaced water acting upward through B. In general, therefore, the prism experiences a torque. In Figure 45B the torque is counterclockwise, and so it turns the prism away from 45A and toward 45C. In 45C the torque vanishes because B is now vertically below C, and this is the orientation that corresponds to stable equilibrium. The torque also vanishes in 45A, and the prism can in principle remain indefinitely in that orientation as well; the equilibrium in this case, however, is unstable, and the slightest disturbance will cause the prism to topple one way or the other. In fact, the potential energy of the system, which increases in a linear fashion with the difference in height between C and B, is at its smallest in orientation 45C and at its largest in orientation 45A. To improve the stability of a floating object one should, if possible, lower C relative to B. In the case of a boat, this may be done by redistributing the load inside.

Archimedes' principle

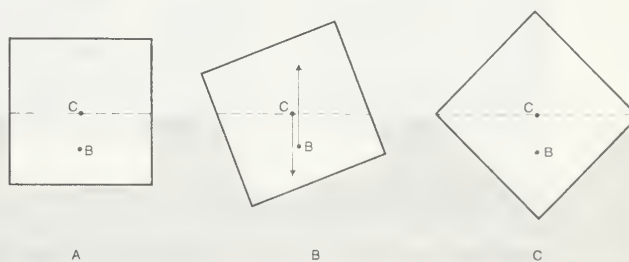


Figure 45: Three possible orientations of a uniform square prism floating in liquid of twice its density. The stable orientation is (C) (see text).

Of the many hydrostatic phenomena in which the surface tension of liquids plays a role, the most significant is probably capillarity. Consider what happens when a tube of narrow bore, often called a capillary tube, is dipped into a liquid. If the liquid "wets" the tube (with zero contact angle), the liquid surface inside the tube forms a concave meniscus, which is a virtually spherical surface having the same radius, r , as the inside of the tube. The tube experiences a downward force of magnitude $2\pi r\sigma$, where σ is the surface tension of the liquid, and the liquid experiences a reaction of equal magnitude that lifts the meniscus through a height h such that

Capillarity

$$2\pi r\sigma = \pi r^2 h \rho g \quad (127)$$

—*i.e.*, until the upward force for which surface tension is responsible is balanced by the weight of the column of liquid that has been lifted. If the liquid does not wet the tube, the meniscus is convex and depressed through the same distance h (see Figure 46). A simple method for determining surface tension involves the measurement of h in one or the other of these situations and the use of equation (127) thereafter.

It follows from equations (124) and (127) that the pressure at a point P just below the meniscus differs from the pressure at Q by an amount

$$\rho g h = \frac{2\sigma}{r}; \quad (128)$$

it is less than the pressure at Q in the case to which Figure 46A refers and greater than the pressure at Q in the other case. Since the pressure at Q is just the atmospheric pressure, it is equal to the pressure at a point immediately above the meniscus. Hence, in both instances there is a pressure difference of $2\sigma/r$ between the two sides of the curved meniscus, and in both the higher pressure is on the inner side of the curve. Such a pressure difference is a requirement of equilibrium wherever a liquid surface is curved. If the surface is curved but not spherical, the pressure difference is

$$\sigma(r_1^{-1} + r_2^{-1}), \tag{129}$$

where r_1 and r_2 are the two principal radii of curvature. If it is cylindrical, one of these radii is infinite, and, if it is curved in opposite directions, then for the purposes of (129) they should be treated as being of opposite sign.

The diagrams in Figure 46 were drawn to represent cross sections through cylindrical tubes, but they might equally well represent two vertical parallel plates that are partly submerged in the liquid a small distance apart. Consideration of how the pressure varies with height shows that over the range of height h the plates experience a greater pressure on their outer surfaces than on their inner surfaces; this is true whether the liquid wets both plates or not. It is a matter of observation that small objects floating near one another on the surface of a liquid tend to move toward one another, and it is the pressure difference just referred to that makes them behave in this way.

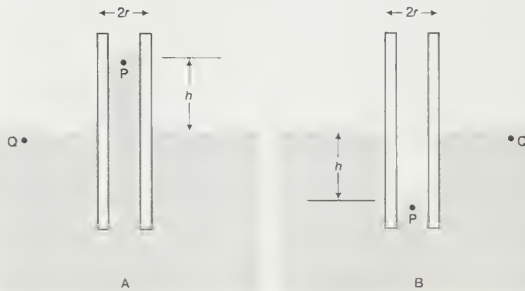


Figure 46: Capillarity. (A) The liquid wets the tube and rises up in it. (B) The liquid does not wet the tube and is depressed.

One other problem having to do with surface tension will be considered here. The diagrams in Figure 47 show stages in the growth of a liquid drop on the end of a tube which the liquid is supposed to wet. In passing from stage A to stage B, by which time the drop is roughly hemispheric in shape, the radius of curvature of the drop diminishes; and it follows from (128) that, to bring about this growth, one must slowly increase the pressure of the liquid inside the tube. If the pressure could be held steady at the value corresponding to B, the drop would then become unstable, because any further growth (e.g., to the more or less spherical shape indicated in Figure 47C) would involve an increase in radius of curvature. The applied pressure would then exceed that required to hold the drop in equilibrium, and the drop would necessarily grow bigger still. In practice, however, it is easier to control the rate of flow of water through the tube, and hence the rate of growth of the drop, than it is to control the pressure. If the rate of flow is very small, drops will form the nonspherical shapes suggested by Figure 47D before they detach themselves and fall. It is not an easy matter to analyze the shape of a drop on the point of detachment, and there is no simple formula for the volume of the drop after it is detached.



Figure 47: Stages in the formation of a liquid drop (see text).

Hydrodynamics

BERNOULLI'S LAW

Up to now the focus has been fluids at rest. This section deals with fluids that are in motion in a steady fashion such that the fluid velocity at each given point in space is not changing with time. Any flow pattern that is steady in this sense may be seen in terms of a set of streamlines, the trajectories of imaginary particles suspended in the fluid

and carried along with it. In steady flow, the fluid is in motion but the streamlines are fixed. Where the streamlines crowd together, the fluid velocity is relatively high; where they open out, the fluid becomes relatively stagnant.

When Euler and Bernoulli were laying the foundations of hydrodynamics, they treated the fluid as an idealized inviscid substance in which, as in a fluid at rest in equilibrium, the shear stresses associated with viscosity are zero and the pressure p is isotropic. They arrived at a simple law relating the variation of p along a streamline to the variation of v (the principle is credited to Bernoulli, but Euler seems to have arrived at it first), which serves to explain many of the phenomena that real fluids in steady motion display. To the inevitable question of when and why it is justifiable to neglect viscosity, there is no single answer. Some answers will be provided later in this article, but other matters will be taken up first.

Consider a small element of fluid of mass m , which—apart from the force on it due to gravity—is acted on only by a pressure p . The latter is isotropic and does not vary with time but may vary from point to point in space. It is a well-known consequence of Newton's laws of motion that, when a particle of mass m moves under the influence of its weight mg and an additional force F from a point P where its speed is v_p and its height is z_p to a point Q where its speed is v_Q and its height is z_Q , the work done by the additional force is equal to the increase in kinetic and potential energy of the particle—*i.e.*, that

$$\int_P^Q F \cdot ds = \left(\frac{1}{2}\right)m(v_Q^2 - v_p^2) + mg(z_Q - z_p). \tag{130}$$

In the case of the fluid element under consideration, F may be related in a simple fashion to the gradient of the pressure, and one finds

$$\int_P^Q F \cdot ds = -m \int_P^Q \rho^{-1} dp. \tag{131}$$

If the variations of fluid density along the streamline from P to Q are negligibly small, the factor ρ^{-1} may be taken outside the integral on the right-hand side of (131), which thereupon reduces to $\rho^{-1}(p_Q - p_p)$. Then (130) and (131) can be combined to obtain

$$\frac{p_p}{\rho} + \frac{v_p^2}{2} + gz_p = \frac{p_Q}{\rho} + \frac{v_Q^2}{2} + gz_Q. \tag{132}$$

Since this applies for any two points that can be visited by a single element of fluid, one can immediately deduce Bernoulli's (or Euler's) important result that along each streamline in the steady flow of an inviscid fluid the quantity

$$\left(\frac{p}{\rho} + \frac{v^2}{2} + gz\right) \tag{133}$$

is constant.

Under what circumstances are variations in the density negligibly small? When they are very small compared with the density itself—*i.e.*, when

$$\left(\frac{\Delta \rho}{\rho}\right) = \beta_s \Delta p = (\beta_s \rho) \Delta \left(\frac{v^2}{2} + gz\right) = \frac{\Delta \left(\frac{v^2}{2} + gz\right)}{V_s^2} \ll 1, \tag{134}$$

where the symbol Δ is used to represent the extent of the change along a streamline of the quantity that follows it, and where V_s is the speed of sound (see below *Compressible flow in gases*). This condition is satisfied for all the flow problems having to do with water that are discussed below. If the fluid is air, it is adequately satisfied provided that the largest excursion in z is on the order of metres rather than kilometres and provided that the fluid velocity is everywhere less than about 100 metres per second.

Bernoulli's law indicates that, if an inviscid fluid is flowing along a pipe of varying cross section, then the pressure is relatively low at constrictions where the velocity is high and relatively high where the pipe opens out and the fluid stagnates. Many people find this situation paradoxical when they first encounter it. Surely, they say, a constriction should increase the local pressure rather than diminish it? The paradox evaporates as one learns to think

of the pressure changes along the pipe as cause and the velocity changes as effect, instead of the other way around; it is only because the pressure falls at a constriction that the pressure gradient upstream of the constriction has the right sign to make the fluid accelerate.

Paradoxical or not, predictions based on Bernoulli's law are well-verified by experiment. Try holding two sheets of paper so that they hang vertically two centimetres or so apart and blow downward so that there is a current of air between them. The sheets will be drawn together by the reduction in pressure associated with this current. Ships are drawn together for much the same reason if they are moving through the water in the same direction at the same speed with a small distance between them. In this case, the current results from the displacement of water by each ship's bow, which has to flow backward to fill the space created as the stern moves forward, and the current between the ships, to which they both contribute, is stronger than the current moving past their outer sides. As another simple experiment, listen to the hissing sound made by a tap that is almost, but not quite, turned off. What happens in this case is that the flow is so constricted and the velocity within the constriction so high that the pressure in the constriction is actually negative. Assisted by the dissolved gases that are normally present, the water cavitates as it passes through, and the noise that is heard is the sound of tiny bubbles collapsing as the water slows down and the pressure rises again on the other side.

Two practical devices that are used by hydraulic engineers to monitor the flow of liquids through pipes are based on Bernoulli's law. One is the venturi tube, a short length with a constriction in it of standard shape (see Figure 48A), which may be inserted into the pipe proper. If the velocity at point P, where the tube has a cross-sectional area A_P , is v_P and the velocity in the constriction, where the area is A_Q , is v_Q , the continuity condition—the condition that the mass flowing through the pipe per unit time has to be the same at all points along its length—suggests that $\rho_P A_P v_P = \rho_Q A_Q v_Q$, or that $A_P v_P = A_Q v_Q$ if the difference between ρ_P and ρ_Q is negligible. Then Bernoulli's law indicates

$$\rho gh = (p_P - p_Q) = \left(\frac{1}{2}\right) \rho v_P^2 \left[\left(\frac{A_P}{A_Q}\right)^2 - 1 \right]. \quad (135)$$

Thus one should be able to find v_P , and hence the quantity $Q (= A_P v_P)$ that engineers refer to as the rate of discharge, by measuring the difference of level h of the fluid in the two side tubes shown in the diagram. At low velocities the pressure difference ($p_P - p_Q$) is greatly affected by viscosity (see below *Viscosity*), and equation (135) is unreliable in

Venturi and pitot tubes

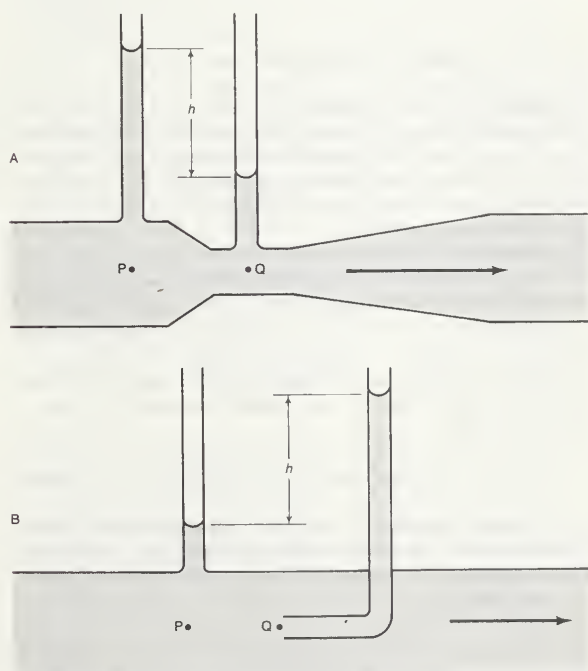


Figure 48: Schematic representation of (A) a venturi tube and of (B) a pitot tube.

consequence. The venturi tube is normally used, however, when the velocity is large enough for the flow to be turbulent (see below *Turbulence*). In such a circumstance, equation (135) predicts values for Q that agree with values measured by more direct means to within a few parts percent, even though the flow pattern is not really steady at all.

The other device is the pitot tube, which is illustrated in Figure 48B. The fluid streamlines divide as they approach the blunt end of this tube, and at the point marked Q in the diagram there is complete stagnation, since the fluid at this point is moving neither up nor down nor to the right. It follows immediately from Bernoulli's law that

$$\rho gh = (p_Q - p_P) = \left(\frac{1}{2}\right) \rho v_P^2. \quad (136)$$

As with the venturi tube, one should therefore be able to find v_P from the level difference h .

One other simple result deserves mention here. It concerns a jet of fluid emerging through a hole in the wall of a vessel filled with liquid under pressure. Observation of jets shows that after emerging they narrow slightly before settling down to a more or less uniform cross section known as the vena contracta. They do so because the streamlines are converging on the hole inside the vessel and are obliged to continue converging for a short while outside. It was Torricelli who first suggested that, if the pressure excess inside the vessel is generated by a head of liquid h , then the velocity v at the vena contracta is the velocity that a free particle would reach on falling through a height h —i.e., that

$$v = \sqrt{(2gh)}. \quad (137)$$

This result is an immediate consequence, for an inviscid fluid, of the principle of energy conservation that Bernoulli's law enshrines.

In the following section, Bernoulli's law is used in an indirect way to establish a formula for the speed at which disturbances travel over the surface of shallow water. The explanation of several interesting phenomena having to do with water waves is buried in this formula. Analogous phenomena dealing with sound waves in gases are discussed below in *Compressible flow in gases*, where an alternative form of Bernoulli's law is introduced. This form of the law is restricted to gases in steady flow but is not restricted to flow velocities that are much less than the speed of sound. The complication that viscosity represents is again ignored throughout these two sections.

WAVES ON SHALLOW WATER

Imagine a layer of water with a flat base that has a small step on its surface, dividing a region in which the depth of the water is uniformly equal to D from a region in which it is uniformly equal to $D(1 + \epsilon)$, with $\epsilon \ll 1$. Let the water in the shallower region flow toward the step with some uniform speed V , as Figure 49A suggests, and let this speed be just sufficient to hold the step in the same position so that the flow pattern is a steady one. The continuity condition (i.e., the condition that as much water flows out to the left per unit time as flows in from the right) indicates that in the deeper region the speed of the water is $V(1 + \epsilon)^{-1}$. Hence by applying Bernoulli's law to the points marked P and Q in the diagram, which lie on the same streamline and at both of which the pressure is atmospheric, one may deduce that

$$geD = \left(\frac{1}{2}\right) V^2 [1 - (1 + \epsilon)^{-2}] \approx \epsilon V^2$$

—i.e., that

$$V = \sqrt{(gD)}. \quad (138)$$

This result shows that, if the water in the shallower region is in fact stationary (see Figure 49B), the step advances over it with the speed V that equation (138) describes, and it reveals incidentally that behind the step the deeper water follows up with speed $V[1 - (1 + \epsilon)^{-1}] \approx \epsilon V$. The argument may readily be extended to disturbances of the surface that are undulatory rather than steplike. Provided that the distance between successive crests—a distance known as the wavelength and denoted by λ —is much

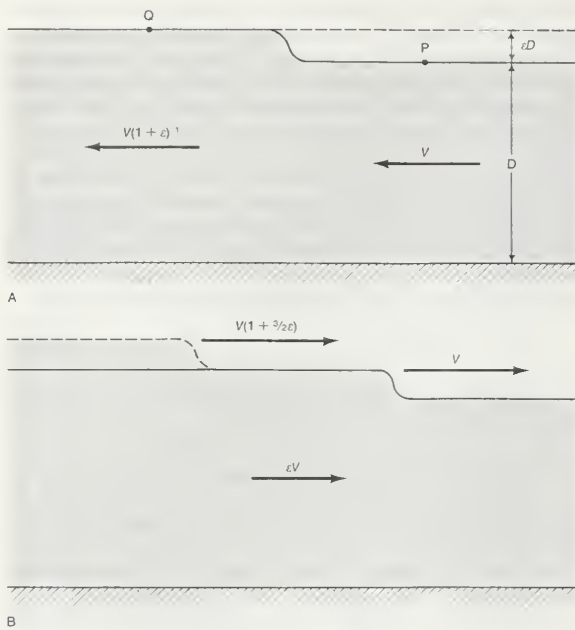


Figure 49: Steps on the surface of shallow water. In (A) the water is moving and the step is stationary. In (B) the water is stationary in front of the first step and the step is therefore moving; the second step (dotted line) is catching up to the first.

Non-dispersive waves

greater than the depth of the water, D , and provided that its amplitude is very much less than D , a wave travels over stationary water at a speed given by (138). Because their speed does not depend on wavelength, the waves are said to be nondispersive.

Evidently waves that are approaching a shelving beach should slow down as D diminishes. If they are approaching it at an angle, the slowing-down effect bends, or refracts, the wave crests so that they are nearly parallel to the shore by the time they ultimately break.

Suppose now that a small step of height ϵD ($\epsilon \ll 1$) is traveling over stationary water of uniform depth D and that behind it is a second step of much the same height traveling in the same direction. Because the second step (suggested by a dotted line in Figure 49B) is traveling on a base that is moving at $\epsilon\sqrt{gD}$ and because the thickness of that base is $(1 + \epsilon)D$ rather than D , the speed of the second step is approximately $(1 + 3\epsilon/2)\sqrt{gD}$. Since this is greater than \sqrt{gD} , the second step is bound to catch up with the first. Hence, if there are a succession of infinitesimal steps that raise the depth continuously from D to some value D' , which differs significantly from D , then the ramp on the surface is bound to become steeper as it advances. It may be shown that if D' exceeds about $1.3D$, the ramp ultimately becomes a vertical step of finite height and that the step then "breaks." A finite step that has broken dissipates energy as heat in the resultant foaming motion, and Bernoulli's equation is no longer applicable to it. A simple argument based on conservation of momentum rather than energy, however, suffices to show that its velocity of propagation is

$$\sqrt{\frac{gD'(D' + D)}{2D}} \tag{139}$$

Tidal bores, which may be observed on some estuaries, are examples on the large scale of the sort of phenomena to which (139) applies. Examples on a smaller scale include the hydraulic jumps that are commonly seen below weirs and sluice gates where a smooth stream of water suddenly rises at a foaming front. In this case, (139) describes the speed of the water, since the front itself is more or less stationary.

When water is shallow but not extremely shallow, so that correction terms of the order of $(D/\lambda)^2$ are significant, waves of small amplitude become slightly dispersive (see below *Waves on deep water*). In this case, a localized disturbance on the surface of a river or canal, which is guided

by the banks in such a way that it can propagate in one direction only, is liable to spread as it propagates. If its amplitude is not small, however, the tendency to spread due to dispersion may in special circumstances be subtly balanced by the factors that cause waves of relatively large amplitude to form bores, and the result is a localized hump in the surface, of symmetrical shape, which does not spread at all. The phenomenon was first observed on a canal near Edinburgh in 1834 by a Scottish engineer named Scott Russell; he later wrote a graphic account of following on horseback, for well over a kilometre, a "large solitary elevation . . . which continued its course along the channel apparently without change of form." What Scott Russell saw is now called a soliton. Solitons on canals can have various widths, but the smaller the width the larger the height must be and the faster the soliton travels. Thus, if a high, narrow soliton is formed behind a low, broad one, it will catch up with the low one. It turns out that, when the high soliton does so, it passes through the low one and emerges with its shape unchanged (see Figure 50).

Solitons

It is now recognized that many of the nonlinear differential equations that appear in diverse branches of physics have solutions of large amplitude corresponding to solitons and that the remarkable capacity of solitons for surviving encounters with other solitons is universal. This discovery has stimulated much interest among mathematicians and physicists, and understanding of solitons is expanding rapidly.

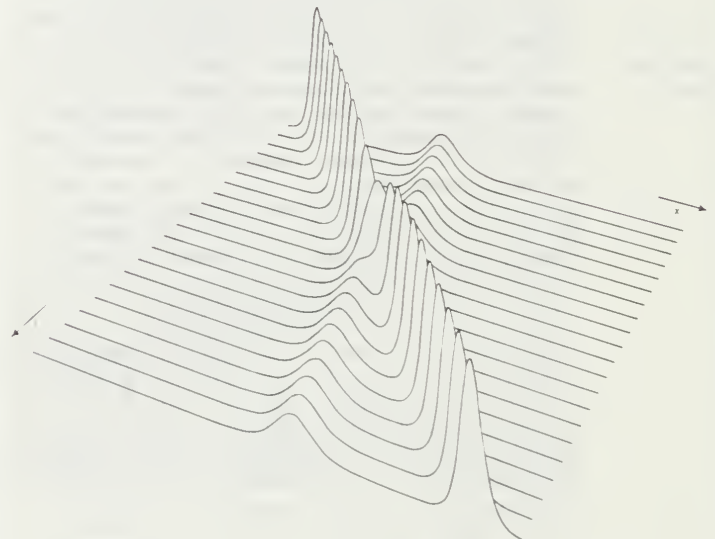


Figure 50: Interaction of two solitons (see text).

COMPRESSIBLE FLOW IN GASES

Compressible flow refers to flow at velocities that are comparable to, or exceed, the speed of sound. The compressibility is relevant because at such velocities the variations in density that occur as the fluid moves from place to place cannot be ignored.

Flow at velocities equal to or greater than the speed of sound

Suppose that the fluid is a gas at a low enough pressure for the ideal equation of state, equation (118), to apply and that its thermal conductivity is so poor that the compressions and rarefactions undergone by each element of the gas may be treated as adiabatic (see above). In this case, it follows from equation (120) that the change of density accompanying any small change in pressure, dp , is such that

$$\rho^{-1} dp = \left(\frac{\gamma}{\gamma - 1} \right) d \left(\frac{p}{\rho} \right) \tag{140}$$

This makes it possible to integrate the right-hand side of equation (131), and one thereby arrives at a version of Bernoulli's law for a steady compressible flow of gases which states that

$$\frac{\gamma p}{(\gamma - 1)\rho} + \frac{v^2}{2} + gz \tag{141}$$

is constant along a streamline. An equivalent statement is that

$$\frac{C_p T}{M} + \frac{v^2}{2} + gz \tag{142}$$

is constant along a streamline. It is worth noting that, when a gas flows through a nozzle or through a shock front (see below), the flow, though adiabatic, may not be reversible in the thermodynamic sense. Thus the entropy of the gas is not necessarily constant in such flow, and as a consequence the application of equation (120) is open to question. Fortunately, the result expressed by (141) or (142) can be established by arguments that do not involve integration of (131). It is valid for steady adiabatic flow whether this is reversible or not.

Bernoulli's law in the form of (142) may be used to estimate the variation of temperature with height in the Earth's atmosphere. Even on the calmest day the atmosphere is normally in motion because convection currents (see below *Convection*) are set up by heat derived from sunlight that is radiated at the Earth's surface. The currents are indeed adiabatic to a good approximation, and their velocity is generally small enough for the term v^2 in (142) to be negligible. One can therefore deduce without more ado that the temperature of the atmosphere should fall off in a linear fashion—i.e., that

$$T(z) = T(0) - \beta z = T(0) - \left(\frac{Mg}{C_p}\right)z. \tag{143}$$

Lapse rate

Here β is used to represent the temperature lapse rate, and the value suggested for this quantity, (Mg/C_p) , is close to 10° C per kilometre for dry air.

This prediction is not exactly fulfilled in practice. Within the troposphere (i.e., to the heights of about 10 kilometres to which convection currents extend), the mean temperature does decrease with height in a linear fashion, but β is only about 6.5° C per kilometre. It is the water vapour in the atmosphere, which condenses as the air rises and cools, that lowers the lapse rate to this value by increasing the effective value of C_p . The fact that the lapse rate is smaller for moist air than for dry air means that a stream of moist air which passes over a mountain range and which deposits its moisture as rain or snow at the summit is warmer when it descends to sea level on the other side of the range than it was when it started. The foehn wind of the Alps owes its warmth to this effect.

The variation of the pressure of the atmosphere with height may be estimated in terms of β , using the equation

$$p(z) = p(0) \left[1 - \frac{\beta z}{T(0)} \right]^{Mg/(\beta R)}. \tag{144}$$

This is obtained by integration of (123), using (118) and (143).

Determining the speed of sound in gases

In the form of equation (141), Bernoulli's law may be used to calculate the speed of sound in gases. The argument is directly analogous to the one applied in the previous section to waves on shallow water—and, indeed, the diagrams in Figure 49 can serve to illustrate the argument here too, if they are regarded as plots of gas density (or else of pressure or temperature, which go hand in hand with density in adiabatic flow) versus position. The results of the argument will be stated without proof. If there exists an infinitesimal step in the density of the gas, it will remain stationary provided that the gas flows uniformly through it toward the region of higher density, with a velocity

$$V_s = \sqrt{\left(\frac{\gamma p}{\rho}\right)}. \tag{145}$$

If the gas is stationary, then (145) describes the velocity with which the step moves. It also describes the speed of propagation of the sort of undulatory variation of density that constitutes a sound wave of fixed frequency or pitch. Because the speed of sound is independent of pitch, sound waves, like waves on shallow water, are nondispersive. This is just as well. It is only because there is no dispersion that one can understand the words of a distant speaker or listen to a symphony orchestra with pleasure from the back of an auditorium as well as from the front.

It should be noted that the formula for the speed of sound in gases may be proved in other ways, and Newton

came close to it a century before Bernoulli's time. However, because Newton failed to appreciate the distinction between adiabatic and isothermal flow, his answer lacked the factor γ occurring in (145). The first person to correct this error was Pierre-Simon Laplace.

The above statements apply to density steps or undulations, the amplitude of which is infinitesimal, and they need some modification if the amplitude is large. In the first place it is found, as for waves on shallow water and for very much the same reasons, that, where two small density steps are moving parallel to one another, the second is bound to catch up with the first. It follows that, if there exists a propagating region in which the density rises in a continuous fashion from ρ to ρ' , where $(\rho' - \rho)$ is not necessarily small, then the width of this region is bound to diminish as time passes. Ultimately a shock front develops over which the density—and hence the pressure and temperature—rises almost discontinuously. There are processes within the shock front, vaguely analogous on the molecular scale to the foaming of a breaking water wave, by which energy is dissipated as heat. The speed of propagation, V_{sh} , of a shock front in a gas that is stationary in front of it may be expressed in terms of V_s and V'_s , the velocities of small-amplitude sound waves in front of the shock and behind it, respectively, by the equation

$$2V_{sh}^2 = \left[\frac{(\gamma + 1)\rho'}{\gamma\rho} \right] V_s'^2 + V_s^2. \tag{146}$$

Thus, if the shock is a strong one ($\rho' \gg \rho$), V_{sh} may be significantly greater than both V_s and V'_s .

Even the gentlest sound wave, in which density and pressure initially oscillate in a smooth and sinusoidal fashion, develops into a succession of weak shock fronts in time. More noticeable shock fronts are a feature of the flow of gases at supersonic speeds through the nozzles of jet engines and accompany projectiles that are moving through stationary air at supersonic speeds. In certain circumstances when a supersonic aircraft is following a curved path, the accompanying shock wave may accidentally reinforce itself in places and thereby become offensively noticeable as a "sonic boom," which may break windowpanes and cause other damage. Strong shock fronts also occur immediately after explosions, of course, and when windowpanes are broken by an explosion, the broken glass tends to fall outward rather than inward. Such is the case because the glass is sucked out by the relatively low density and pressure that succeed the shock itself.

The diagrams in Figure 51 show a well-known construction attributed to the Austrian physicist Ernst Mach that explains the origin of the shock front accompanying a supersonic projectile. The circular arcs in this figure represent cross sections through spherical disturbances that are spreading with speed V_s from centres (S' , S'' , etc.), which mark the position of the moving source S at the time when they were emitted. If the source is something like the tip of

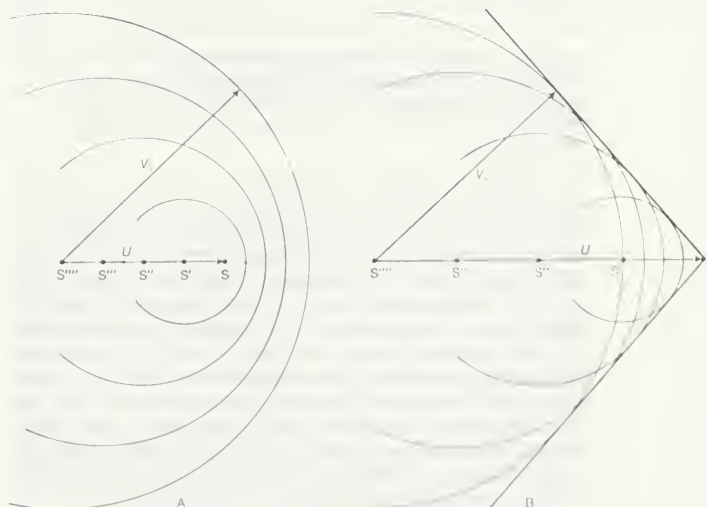


Figure 51: Mach's construction. (A) Source speed U less than speed of sound V_s , (B) U greater than V_s (see text).

an arrow, which disturbs the air by parting it as it travels along but which is inaudible when stationary, then each "disturbance" due to some infinitesimal displacement of the tip is a spherical shell of infinitesimal thickness within which a small radial velocity has been imparted to the air. There is an infinite number of such disturbances, overlapping one another, of which only a handful are represented in Figure 51. When the velocity of the source, U , is less than V_s (Figure 51A), the result of adding them together is the sort of steady backflow that is to be expected around a moving obstacle, and there is no sound emission in the normal sense; the source remains inaudible. When U exceeds V_s , however, the spherical disturbances reinforce one another, as Figure 51B shows, on a conical caustic surface, which makes an angle of $\sin^{-1}(U/V)$ to the line of travel of the source, and it is on this surface that a shock front is to be expected. The cone becomes sharper as the source speeds up.

VISCOSITY

As shown above, a number of phenomena of considerable physical interest can be discussed using little more than the law of conservation of energy, as expressed by Bernoulli's law. However, the argument has so far been restricted to cases of steady flow. To discuss cases in which the flow is not steady, an equation of motion for fluids is needed, and one cannot write down a realistic equation of motion without facing up to the problems presented by viscosity, which have so far been deliberately set aside.

The concept of viscosity was first formalized by Newton, who considered the shear stresses likely to arise when a fluid undergoes what is called laminar motion with the sort of velocity profile that is suggested in Figure 52A; the laminae here are planes normal to the x_2 -axis, and they are moving in the direction of the x_1 -axis with a velocity v_1 , which increases in a linear fashion with x_2 . Newton suggested that, as each lamina slips over the one below, it exerts a sort of frictional force upon the latter in the forward direction, in which case the upper lamina is bound to experience an equal reaction in the backward direction. The strength of these forces per unit area constitutes the

Stresses in laminar motion

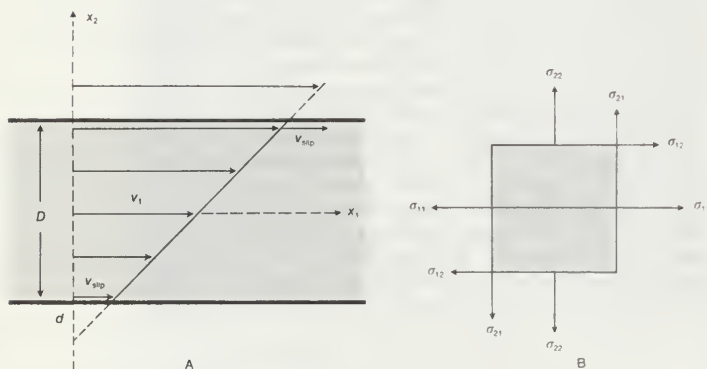


Figure 52: Laminar motion and associated stresses

(A) Velocity profile for laminar flow between two plates, driven by motion of the upper plate, (B) an enlarged view of a cubic element of the fluid between the plates, showing the stresses that act upon it.

component of shear stress normally written as σ_{12} (not to be confused with surface tension, for which the symbol σ has been used above). Figure 52B shows, in elevation, an enlarged view of an infinitesimal element of the fluid of cubic shape, and the directions of the forces experienced by this cube associated with σ_{12} are indicated by arrows. Other arrows show the directions of the forces associated with the so-called normal stresses σ_{11} and σ_{22} , which in the absence of motion of the fluid would both be equal, by Pascal's law, to $-p$. Now σ_{12} is clearly zero when the rate of variation of velocity, $\partial v_1/\partial x_2$, is zero, for then there is no slip, and presumably it increases monotonically as $\partial v_1/\partial x_2$ increases. Newton made the plausible assumption that the two are linearly related—*i.e.*, that

$$\sigma_{12} = \eta \left(\frac{\partial v_1}{\partial x_2} \right). \tag{147}$$

The full name for the coefficient η is shear viscosity to distinguish it from the bulk viscosity, b , which is defined below. The word shear, however, is frequently omitted in this context.

Now if the only shear stress acting on the cubic element of fluid sketched in Figure 52B were σ_{12} , the cube would experience a torque tending to make it twist in a clockwise sense. Since the magnitude of the torque would vary like the third power of the linear dimensions of the cube, whereas the moment of inertia of the element would vary like the fifth power, the resultant angular acceleration for an infinitesimal cube would be infinite. One may infer that any tendency to twist in a clockwise sense gives rise instantaneously to an additional shear stress σ_{21} , the direction of which is indicated in the diagram, and that σ_{12} and σ_{21} are equal at all times. It follows that equation (147) cannot be a complete expression for these shear stresses, for it does not include the possibility that the fluid is moving in the x_2 direction, with a velocity v_2 that varies with x_1 . The complete expression for what is called a Newtonian fluid is

$$\sigma_{12} = \sigma_{21} = \eta \left[\left(\frac{\partial v_1}{\partial x_2} \right) + \left(\frac{\partial v_2}{\partial x_1} \right) \right]. \tag{148}$$

Similar expressions may be written down for σ_{23} ($=\sigma_{32}$) and σ_{31} ($=\sigma_{13}$). Since Newton's day these hypothetical expressions have been fully substantiated for gases and simple liquids, not only by experiment but also by analysis of the molecular motions and molecular interactions in such fluids undergoing shear, and for such fluids one can even predict the magnitude of η with reasonable success. There do exist, however, more complicated fluids for which the Newtonian description of shear stress is inadequate, and some of these are very familiar in the home. In the whites of eggs, for example, and in most shampoos, there are long-chain molecules that become entangled with one another, and entanglement may hinder their efforts to respond to changes of environment associated with flow. As a result, the stresses acting in such fluids may reflect the deformations experienced by the fluid in the recent past as much as the instantaneous rate of deformation. Moreover, the relation between stress and rate of deformation may be far from linear. Non-Newtonian effects, interesting though they are, lie outside the scope of the present discussion, however.

Non-Newtonian effects

The sort of velocity profile that is suggested by Figure 52B may be established by containing the fluid between two parallel flat plates and moving one plate relative to the other. The possibility exists that in this situation the layers of fluid immediately in contact with each plate will slip over them with some finite velocity (indicated in the diagram by an arrow labeled v_{slip}). If so, the frictional stresses associated with this slip must be such as to balance the shear stress $\eta(\partial v_1/\partial x_2)$ exerted on each of these layers by the rest of the fluid. Little is known about fluid-solid frictional stresses, but intelligent guesswork suggests that they are proportional in magnitude to v_{slip} and that, in the circumstances to which Figure 52A refers, the distance d below the surface of the stationary bottom plate at which the straight line representing the variation of v_1 with x_2 extrapolates to zero should be of the same order of magnitude as the diameter of a molecule if the fluid is a liquid or as the molecular "mean free path" if it is a gas. These distances are normally very small compared with the separation of the plates, D . Accordingly, fluid flow patterns may normally be treated as subject to the boundary condition that at a fluid-solid interface the relative velocity of the fluid is zero. No reliable evidence for failure of predictions based on this no-slip boundary condition has yet been found, except in the case of what is called Knudsen flow of gases (*i.e.*, flow at such low pressures that the mean free path is comparable in length with the dimensions of the apparatus).

If a fluid is flowing steadily between two parallel plates that are both stationary and if its velocity must be zero in contact with both of them, the velocity profile must necessarily have the form indicated in Figure 53. A force in the forward direction due to the shear stress $\eta(\partial v_1/\partial x_2)$ is transmitted to the plates, and an equal force in the

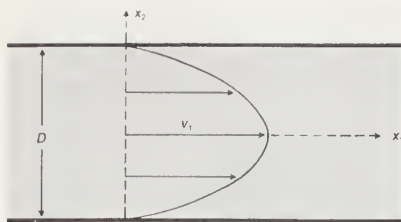


Figure 53: Velocity profile for laminar flow between two plates (or inside a cylindrical tube), driven by a pressure gradient (see text).

backward direction acts on the fluid. The motion therefore cannot be maintained unless the pressure acting on the fluid is greater on the left of the diagram than it is on the right. A full analysis shows the velocity profile to be parabolic, and it indicates that the rate of discharge is related to the pressure gradient by the equation

$$Q = -\left(\frac{WD^3}{12\eta}\right)\left(\frac{dp}{dx_1}\right), \tag{149}$$

where W ($\gg D$) is the width of the plates, measured perpendicular to the diagram in Figure 53. A similar analysis of the problem of steady flow through a (horizontal) cylindrical pipe of uniform diameter D , to which Figure 53 could equally well apply, shows the rate of discharge in this case to be given by

$$Q = -\left(\frac{\pi D^4}{128\eta}\right)\left(\frac{dp}{dx_1}\right), \tag{150}$$

Poiseuille's equation

this famous result is known as Poiseuille's equation, and the type of flow to which it refers is called Poiseuille flow.

Viscosity may affect the normal stress components, σ_{11} , σ_{22} , and σ_{33} , as well as the shear stress components. To see why this is so, one needs to examine the way in which stress components transform when one's reference axes are rotated. Here, the result will be stated without proof that the general expression for σ_{11} consistent with (148) is

$$\begin{aligned} \sigma_{11} = & -p + \left(b + \frac{4\eta}{3}\right)\left(\frac{\partial v_1}{\partial x_1}\right) + \left(b - \frac{2\eta}{3}\right)\left(\frac{\partial v_2}{\partial x_2}\right) \\ & + \left(b - \frac{2\eta}{3}\right)\left(\frac{\partial v_3}{\partial x_3}\right). \end{aligned} \tag{151}$$

On the right-hand side of this equation, p represents the equilibrium pressure defined in terms of local density and temperature by the equation of state, and b is another viscosity coefficient known as the bulk viscosity.

The bulk viscosity is relevant only where the density is changing. Thus it plays a role in attenuating sound waves in fluids and may be estimated from the magnitude of the attenuation. If the fluid is effectively incompressible, however, so that changes of density may be ignored, the flow is everywhere subject to the continuity condition that

$$\left(\frac{\partial v_1}{\partial x_1}\right) + \left(\frac{\partial v_2}{\partial x_2}\right) + \left(\frac{\partial v_3}{\partial x_3}\right) \equiv \nabla \cdot \mathbf{v} \text{ or } \text{div } \mathbf{v} = 0. \tag{152}$$

The terms in (151) that involve b then cancel, and the expression simplifies to

$$\sigma_{11} = -p + 2\eta\left(\frac{\partial v_1}{\partial x_1}\right). \tag{153}$$

Similar equations may be written down for σ_{22} and σ_{33} . These simpler expressions provide the basis for the argument that follows, and the bulk viscosity can be left on one side.

A variety of methods are available for the measurement of shear viscosity. One standard method involves measurement of the pressure gradient along a pipe for various rates of flow and application of Poiseuille's equation. Other methods involve measurement either of the damping of the torsional oscillations of a solid disk supported between two parallel plates when fluid is admitted to the space between the plates, or of the effect of the fluid on the frequency of the oscillations.

The Couette viscometer deserves a fuller explanation. In this device, the fluid occupies the space between two coaxial cylinders of radii a and b ($b > a$); the outer cylinder is

rotated with uniform angular velocity ω_0 , and the resultant torque transmitted to the inner stationary cylinder is measured. If both the terms on the right-hand side of equation (148) are taken into account, the shear stress in the circulating fluid is found to be proportional to $r(d\omega/dr)$ rather than to (dv/dr) —not an unexpected result, since it is only if ω , the angular velocity of the fluid, varies with radius r that there is any slip between one cylindrical lamina of fluid and the next. The torque transmitted through the fluid is therefore proportional to $r^3(d\omega/dr)$. In the steady state, the opposing torques acting on the inner and outer surfaces of each cylindrical lamina of fluid must be of equal magnitude—otherwise the laminae accelerate—and this means that $r^3(d\omega/dr)$ must be independent of r . There are two basic modes of motion for a circulating fluid that satisfy this condition: in one, the liquid rotates as a solid body would, with an angular velocity that does not vary with r , and the torque is everywhere zero; in the other, ω varies like r^{-2} . The angular velocity of the fluid in a Couette viscometer can be viewed as a mixture of these two modes in proportions that satisfy the boundary conditions at $r = a$ and $r = b$. The torque transmitted per unit length of the cylinders turns out to be given by

$$4\pi\eta\omega_0\left[\frac{b^2a^2}{(b^2 - a^2)}\right]. \tag{154}$$

It may be added that if the inner cylinder is absent, the steady flow pattern consists only of the first mode—*i.e.*, the fluid rotates like a solid body with uniform angular velocity ω_0 . If the outer cylinder is absent, however, and the inner one rotates, it then consists only of the second mode. The angular velocity falls off like r^{-2} , and the velocity v falls off like r^{-1} .

In the equation of motion given in the following section, the shear viscosity occurs only in the combination (η/ρ) . This combination occurs so frequently in arguments of fluid dynamics that it has been given a special name—kinetic viscosity. The kinetic viscosity at normal temperatures and pressures is about 10^{-6} square metre per second for water and about 1.5×10^{-5} square metre per second for air.

Kinetic viscosity

NAVIER-STOKES EQUATION

One may have a situation where σ_{11} increases with x_1 . The force that this component of stress exerts on the right-hand side of the cubic element of fluid sketched in Figure 52B will then be greater than the force in the opposite direction that it exerts on the left-hand side, and the difference between the two will cause the fluid to accelerate along x_1 . Accelerations along x_1 will also result if σ_{12} and σ_{13} increase with x_2 and x_3 , respectively. These accelerations, and corresponding accelerations in the other two directions, are described by the equation of motion of the fluid. For a fluid moving so slowly compared with the speed of sound that it may be treated as incompressible and in which the variations of temperature from place to place are insufficient to cause significant variations in the shear viscosity η , this equation takes the form

$$\begin{aligned} -\nabla\left(\frac{p}{\rho} + gz\right) - \left(\frac{\eta}{\rho}\right)[\nabla \times (\nabla \times \mathbf{v})] &= \frac{D\mathbf{v}}{Dt} \\ &= \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla)\mathbf{v}. \end{aligned} \tag{155}$$

Euler derived all the terms in this equation except the one on the left-hand side proportional to (η/ρ) , and without that term the equation is known as the Euler equation. The whole is called the Navier-Stokes equation.

The equation is written in a compact vector notation which many readers will find totally impenetrable, but a few words of explanation may help some others. The symbol ∇ represents the gradient operator, which, when preceding a scalar quantity X , generates a vector with components $(\partial X/\partial x_1, \partial X/\partial x_2, \partial X/\partial x_3)$. The vector product of this operator and the fluid velocity \mathbf{v} —*i.e.*, $(\nabla \times \mathbf{v})$ —is sometimes designated as **curl** \mathbf{v} [and $\nabla \times (\nabla \times \mathbf{v})$ is also **curl curl** \mathbf{v}]. Another name for $(\nabla \times \mathbf{v})$, which expresses particularly vividly the characteristics of the local flow pattern that it represents, is vorticity. In a sample of fluid

Vorticity

that is rotating like a solid body with uniform angular velocity ω_0 , the vorticity lies in the same direction as the axis of rotation, and its magnitude is equal to $2\omega_0$. In other circumstances the vorticity is related in a similar fashion to the local angular velocity and may vary from place to place. As for the right-hand side of (155), Dv/Dt represents the rate of change of velocity that one would see if the motion of a single element of the fluid could be followed—that is, it represents the acceleration of the element—while $\partial v/\partial t$ represents the rate of change at a fixed point in space. If the flow is steady, then $\partial v/\partial t$ is everywhere zero, but the fluid may be accelerating all the same, as individual fluid elements move from regions where the streamlines are widely spaced to regions where they are close together. It is the difference between Dv/Dt and $\partial v/\partial t$ —i.e., the final $(v \cdot \nabla)v$ term in (155)—that introduces into fluid dynamics the nonlinearity that makes the subject so rife with surprises.

POTENTIAL FLOW

This section is concerned with an important class of flow problems in which the vorticity is everywhere zero, and for such problems the Navier-Stokes equation may be greatly simplified. For one thing, the viscosity term drops out of it. For another, the nonlinear term, $(v \cdot \nabla)v$, may be transformed into $\nabla(v^2/2)$. Finally, it may be shown that, when $(\nabla \times v)$ is zero, one may describe the velocity by means of a scalar potential ϕ , using the equation

$$v = \nabla\phi \quad [\equiv \text{grad } \phi]. \quad (156)$$

Thus (155) becomes

$$-\nabla\left(\frac{p}{\rho} + gz + \frac{v^2}{2} + \frac{\partial\phi}{\partial t}\right) = 0,$$

which may at once be integrated to show that

$$\left(\frac{p}{\rho} + gz + \frac{v^2}{2} + \frac{\partial\phi}{\partial t}\right) = \text{constant}. \quad (157)$$

This result incorporates Bernoulli's law for an effectively incompressible fluid ([133]), as was to be expected from the disappearance of the viscosity term. It is more powerful than (133), however, because it can be applied to nonsteady flow in which $\partial\phi/\partial t$ is not zero and because it shows that in cases of potential flow the left-hand side of (157) is constant everywhere and not just constant along each streamline.

Vorticity-free, or potential, flow would be of rather limited interest were it not for the theorem, first proved by Thomson, that, in a body of fluid which is free of vorticity initially, the vorticity remains zero as the fluid moves. This theorem seems to open the door for relatively painless solutions to a great range of problems. Consider, for example, a stream of fluid in uniform motion approaching an obstacle of some sort. Well upstream of the obstacle the fluid is certainly vorticity-free, so it should, according to Thomson's theorem, be vorticity-free around the obstacle and downstream as well. In this case a flow potential should exist; and, if the fluid is effectively incompressible, it follows from equations (152) and (156) that it satisfies Laplace's equation,

$$\left(\frac{\partial^2\phi}{\partial x_1^2}\right) + \left(\frac{\partial^2\phi}{\partial x_2^2}\right) + \left(\frac{\partial^2\phi}{\partial x_3^2}\right) [\equiv \nabla^2\phi] = 0. \quad (158)$$

This is perhaps the most frequently occurring differential equation in physics, and methods for solving it, subject to appropriate boundary conditions, are very well established. Given a solution for ϕ , the fluid velocity v follows at once, and one may then discover how the pressure varies with position and time from equation (157).

The physicists and mathematicians who developed fluid dynamics during the 19th century relied heavily on this reasoning. They based splendid achievements upon it, a notable example being the theory of waves on deep water (see below). There was a touch of unreality, however, about some of their theorizing. If carried to extremes, the argument of the previous section implies that water initially stationary in a beaker can never be set into rotation by rotating the beaker or by stirring it with a spoon, and this is clearly nonsense. It suggests that vorticity-free

water remains vorticity-free if it is squeezed into a narrow pipe, and this too is plainly nonsensical, for the well-established parabolic profile illustrated by Figure 53 is not vorticity-free. What is misleading about the argument in situations like these is that it pays inadequate attention to what happens at interfaces. Following the work of Prandtl, physicists now appreciate that vorticity is liable to be fed into the fluid at interfaces, whether these are interfaces between the fluid and some solid object or the free surfaces of a liquid. Once the slightest trace of vorticity is present, it destroys the conditions on which the proof of Thomson's theorem depends. Moreover, vorticity admitted at interfaces spreads into the fluid in much the same way that a dye would spread, and whether or not the results of potential theory are useful depends on how much of the fluid is contaminated in the particular circumstances under discussion.

POTENTIAL FLOW WITH CIRCULATION: VORTEX LINES

The proof of Thomson's theorem depends on the concept of circulation, which Thomson introduced. This quantity is defined for a closed loop which is embedded in, and moves with, the fluid; denoted by K , it is the integral around the loop of $v \cdot dl$, where dl is an element of length along the loop. If the vorticity is everywhere zero, then so is the circulation around all possible loops, and vice versa. Thomson showed that K cannot change if the viscous term in (155) contributes nothing to the local acceleration, and it follows that both K and vorticity remain zero for all time.

Reference was made earlier to the sort of steady flow pattern that may be set up by rotating a cylindrical spindle in a fluid; the streamlines are circles around the spindle, and the velocity falls off like r^{-1} . This pattern of flow occurs naturally in whirlpools and typhoons, where the role of the spindle is played by a "core" in which the fluid rotates like a solid body; the axis around which the fluid circulates is then referred to as a vortex line. Each small element of fluid outside the core, if examined in isolation for a short interval of time, appears to be undergoing translation without rotation, and the local vorticity is zero. Were it not so, the viscous torques would not cancel and the flow pattern would not be a steady one. Nevertheless, the circulation is not zero if the loop for which it is defined is one that encloses the spindle or core. In such situations, a potential that obeys Laplace's equation outside the spindle or core can be found, but it is no longer, to use a technical term that may be familiar to some readers, single-valued.

Readers who recognize this term are likely to have encountered it in the context of electromagnetism, and it is worth remarking that all the results of potential flow theory have electromagnetic analogues, in which streamlines become the lines of force of a magnetic field and vortex lines become lines of electric current. The analogy may be illustrated by reference to the Magnus effect.

This effect (named for the German physicist and chemist H.G. Magnus, who first investigated it experimentally) arises when fluid flows steadily past a cylindrical spindle, with a velocity that at large distances from the spindle is perpendicular to the spindle's axis and uniformly equal to, say, v_0 , while the spindle itself is steadily rotated. Rotation is communicated to the fluid, and in the steady state the circulation around any loop that encloses the spindle (and encloses a layer of fluid adjacent to the spindle within which the vorticity is nonzero and potential theory is inapplicable) has some nonzero value K . The streamlines that describe the steady flow pattern (outside that "boundary layer") have the form suggested by Figure 54, though the details naturally depend on the magnitude of v_0 and K . The flow pattern has stagnation points at P and P' and, since the pressure is high at such points, the spindle may be expected to experience a downward force perpendicular both to its axis and to the direction of v_0 . Detailed calculations confirm this expectation and show that the magnitude of the force, per unit length of the spindle, is

$$\rho v_0 K \quad (159)$$

This so-called Magnus force is directly analogous to the force that a transverse magnetic field B_0 exerts upon a wire

Vortex lines

Magnus effect

Thomson's theorem

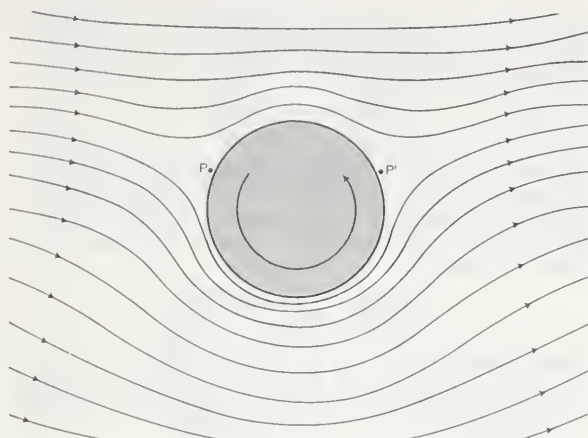


Figure 54: Streamlines for potential flow with circulation past a rotating cylinder. The cylinder experiences a downward Magnus force (see text).

carrying an electric current I , the magnitude of which, per unit length of the wire, is $B_0 I$.

The Magnus force on rotating cylinders has been utilized to propel experimental yachts, and it is closely related to the lift force on airfoils that enables airplanes to fly (see below *Lift*). The transverse forces that cause spinning balls to swerve in flight are, however, not Magnus forces, as is sometimes asserted. They are due to the asymmetrical nature of the eddies that develop at the rear of a spinning sphere (see below *Boundary layers and separation*). Cricket balls, unlike the balls used for baseball, tennis, and golf, have a raised equatorial seam that plays an important part in making the eddies asymmetric. A bowler in cricket who wants to make the ball swerve imparts spin to it, but he does so chiefly to ensure that the orientation of this seam remains steady as the ball moves toward the batsman.

It may be shown, by reference to the magnetic analogue or in other ways, that straight vortex lines of equal but opposite strength, $\pm K$, which are parallel and separated by a distance d , will drift sideways together through the fluid at a speed given by $K/2\pi d$. Similarly, a vortex line that has joined up on itself to form a closed vortex ring of radius a drifts along its axis with a speed given by

$$\left(\frac{K}{4\pi a}\right)\ln\left(\frac{a}{c}\right), \quad (160)$$

where c is the radius of the line's core, with \ln standing for natural logarithm. This formula applies, for example, to smoke rings. The fact that such rings slow down as they propagate can be explained in terms of the increase of c with time, due to viscosity.

WAVES ON DEEP WATER

One particular solution of Laplace's equation that describes wave motion on the surface of a lake or of the ocean is

$$\varphi = \varphi_0 \cos\left\{2\pi\left[\left(\frac{x}{\lambda}\right) - ft\right]\right\} \sinh\left[2\pi\frac{(D+z)}{\lambda}\right]. \quad (161)$$

In this case the x -axis is the direction of propagation and the z -axis is vertical; $z=0$ describes the free surface of the water when it is undisturbed and $z=-D$ describes the bottom surface; φ_0 is an arbitrary constant that determines the amplitude of the motion; and f is the frequency of the waves and λ their wavelength. If λ is more than a few centimetres, surface tension is irrelevant and the pressure in the liquid just below its free surface is atmospheric for all values of x . It can be shown that in these circumstances the wave motion described by (161) is consistent with (157) only if the frequency and wavelength are related by the equation

$$f^2 = \left(\frac{2\pi g}{\lambda}\right) \tanh\left(\frac{2\pi D}{\lambda}\right), \quad (162)$$

and an expression for the speed of the waves may be deduced from this, since $V=f\lambda$. For shallow water ($D \ll \lambda$) one obtains the answer already quoted as equation (138), but for deep water ($D \gg \lambda$) the answer is

$$V = \sqrt{\left(\frac{g\lambda}{2\pi}\right)}. \quad (163)$$

Waves on deep water are evidently dispersive, and surfers rely on this fact. A storm in the middle of the ocean disturbs the surface in a chaotic way that would be useless for surfing, but as the component waves travel toward the shore they separate: those with long wavelengths move ahead of those with short wavelengths because they travel faster. As a result, the waves seem nicely regular by the time that they arrive.

Dispersive character of waves on deep water

Anyone who has observed the waves behind a moving ship will know that they are confined to a V-shaped area of the water's surface, with the ship at its apex. The waves are particularly prominent on the arms of the V, but they can also be discerned between these arms where the wave crests curve in the manner indicated in Figure 55. It seems to be widely believed that the angle of the V becomes more acute as the boat speeds up, much in the way that the conical shock wave accompanying a supersonic projectile becomes more acute (see Figure 51). That is not the case; the dispersive character of waves on deep water is such that the V has a fixed angle of $2 \sin^{-1}(1/3) = 39^\circ$. Thomson (Lord Kelvin) was the first to explain this, and so the V-shaped area is now known as the Kelvin wedge.

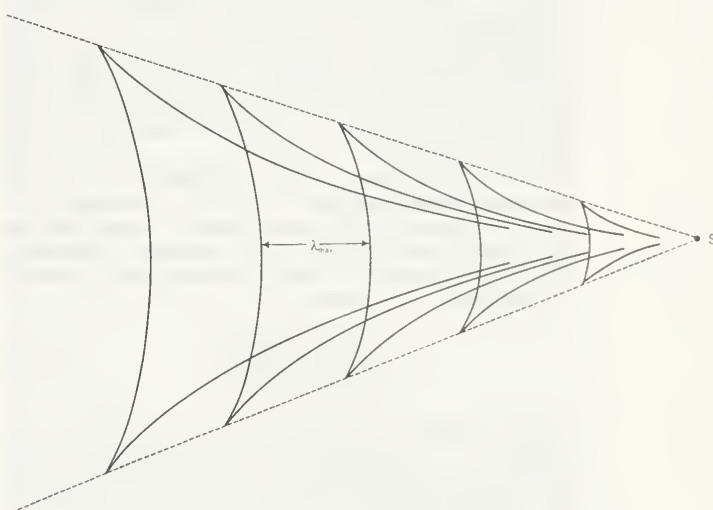


Figure 55: Wave crests in the Kelvin wedge behind a source S that is moving steadily from left to right. The maximum wavelength λ_{\max} depends on the speed of the source, but the angle of the wedge does not (see text).

A version of Thomson's argument is illustrated by the diagram in Figure 56. Here S (the "source") represents the bow of the ship which is moving from left to right with uniform speed U , and the lines labeled $C, C', C'',$ etc., represent a set of parallel wave crests which are also moving from left to right. It can be shown that S will create this set of crests if, but only if, it rides continuously on the one labeled C . (It also can be shown that, though the crests in the set continue indefinitely to the left of C , there can be none to the right of this one.) The condition that S and C move together indicates that there is a relation between wavelength λ and inclination α expressed by the equation

$$\sin \alpha = \frac{V}{U} = \sqrt{\frac{g\lambda}{2\pi U^2}}. \quad (164)$$

This condition can evidently be satisfied by many other sets of crests besides the one represented by full lines in the figure—e.g., by the set with slightly shorter wavelength λ' that is represented by broken lines. When one takes into consideration all the sets that satisfy (164) and have wavelengths intermediate between λ and λ' , it becomes apparent that over most of the area behind the source they interfere destructively. They reinforce one another, however, near the intersections that are ringed in the figure. These intersections lie on a line through S of inclination β , where

$$\tan \beta = \frac{\tan \alpha}{(2 + \tan^2 \alpha)}. \quad (165)$$

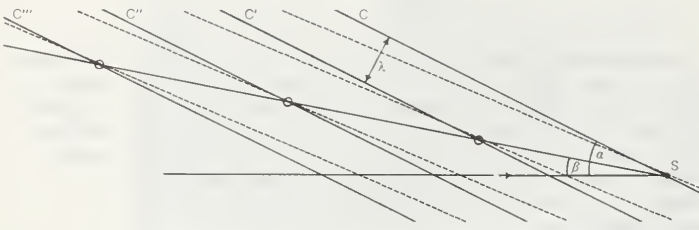


Figure 56: The curved wave crests of Figure 55 result from the superposition of many sets of straight wave crests like the two shown here. These two sets and others that are intermediate in wavelength reinforce one another near the line of inclination β and interfere destructively elsewhere.

It follows that, though the angle α can take any value between 90° (corresponding to $\lambda = \lambda_{\max} = 2\pi U^2/g$) and zero, $\tan \beta$ can never exceed $1/2\sqrt{2}$, and $\sin \beta$ can never exceed $1/3$.

Ships lose energy to the waves in the Kelvin wedge, and they experience additional resistance on that account. The resistance is particularly high when the wave system created by the bow, where water is pushed aside, reinforces the wave system created by the "anti-source" at the stern, where the water closes in again. Such reinforcement is liable to occur when the effective length of the boat, L , is equal to $(2n+1)\lambda_{\max}/2$ (with $n=0, 1, 2, \dots$) and therefore when the Froude number, U/\sqrt{Lg} , takes one of the values $[\sqrt{(2n+1)\pi}]^{-1}$. However, once a boat has been accelerated past $U = \sqrt{Lg/\pi}$, the bow and stern waves tend to cancel, and the resistance resulting from wave creation diminishes.

Waves on deep water whose wavelength is a few centimetres or less are generally referred to as ripples. In such waves, the pressure differences across the curved surface of the water associated with surface tension (see equation [129]) are not negligible, and the appropriate expression for their speed of propagation is

$$V = \sqrt{\left(\frac{g\lambda}{2\pi}\right) + \left(\frac{2\pi\sigma}{\lambda\rho}\right)}. \quad (166)$$

The wave velocity is therefore large for very short wavelengths as well as for very long ones. For water at normal temperatures, V has a minimum value of about 0.23 metre per second when the wavelength is about 17 millimetres, and it follows (note that equation [164] has no real root for α unless U exceeds V) that an object moving through water can create no ripples at all unless its speed exceeds 0.23 metre per second. A wind moving over the surface of water likewise creates no ripples unless its speed exceeds a certain critical value, but this is a more complicated phenomenon, and the critical speed in question is distinctly higher.

BOUNDARY LAYERS AND SEPARATION

It should be reiterated that vorticity is liable to enter a fluid that is initially undergoing potential flow where it makes contact with a solid and also at its free surface. The way in which, having entered, it spreads, may be illustrated by a simple example. Consider a large body of fluid, initially stationary, being set into motion by the movement in its own plane of a large solid plate that is immersed within the fluid. The motion is communicated from solid to fluid by the frictional forces that prevent slip between the two (see above *Viscosity*), and a velocity

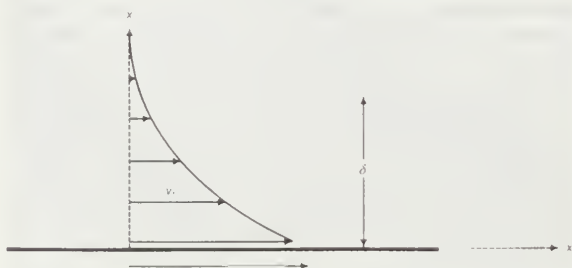


Figure 57: Velocity profile established by motion of a plate through stationary fluid (see text).

profile of the form suggested by Figure 57 is established. Its development with time turns out to be described by the partial differential equation

$$\rho \left(\frac{\partial v_1}{\partial t} \right) = \eta \left(\frac{\partial^2 v_1}{\partial x_2^2} \right). \quad (167)$$

In this situation the vorticity, which may be denoted by the symbol Ω , has one nonzero component, directed along the axis perpendicular to the diagram in Figure 57: it is $\Omega_3 = -(\partial v_1/\partial x_2)$. Differentiation of (167) with respect to x_2 shows at once that

$$\rho \left(\frac{\partial \Omega_3}{\partial t} \right) = \eta \left(\frac{\partial^2 \Omega_3}{\partial x_2^2} \right). \quad (168)$$

This is a diffusion equation. It indicates that, if the plate oscillates to and fro with frequency f , then the so-called boundary layer within which Ω_3 is nonzero has a thickness δ given by

$$\delta \approx \sqrt{\left(\frac{\eta}{\pi\rho f}\right)}. \quad (169)$$

and in most instances of oscillatory motion this is small enough for the boundary layer to be neglected. For example, the boundary layer on the surface of the ocean has a thickness of less than one millimetre when a wave with a frequency of about one hertz passes by: because the effects of viscosity are confined to this layer, they are too slight to affect the propagation of the wave to any significant degree. If the plate is kept moving at a uniform rate, however, the thickness of the boundary layer, as described by (168), will increase with the time t that has elapsed since the motion of the plate began, according to the equation

$$\delta \approx \sqrt{\left(\frac{2\eta t}{\rho}\right)}. \quad (170)$$

Prandtl suggested that when a stream of fluid flows steadily past an obstacle of finite extent, such as a sphere, the time that matters is the time for which fluid on a streamline just outside the boundary layer remains in contact with it. This time is of order D/v_0 , where D is the diameter of the sphere and v_0 is the speed of the fluid well upstream. Hence, one would expect the thickness of the boundary layer at the rear of the sphere to be something like

$$\sqrt{\left(\frac{\eta D}{\rho v_0}\right)}. \quad (171)$$

If the velocity v_0 is so low that (170) is comparable with or greater than the diameter D , the flow pattern must be so contaminated by vorticity that the neglect of viscosity and reliance on Bernoulli's equation and on the other results of potential theory is clearly unjustified. If the velocity is high and (171) is much less than D , however, the boundary layer would seem to be of little importance. Surely then the results of potential theory are to be trusted?

Alas, that optimistic conclusion is not confirmed by experiment. What happens at high velocities is that the boundary layer comes unstuck from the surface of the sphere—it is said to separate. The reason why it does so is suggested by Figure 58A, which shows the streamlines to be expected when the boundary layer (shown in this figure by a shaded area still attached to the sphere) is relatively thin. Evidently the fluid velocity is higher near the equator of the sphere, at Q, than it is at either of the two poles, P and P'. Thus according to Bernoulli's equation, which can be relied on outside the boundary layer, the pressure near Q is less than it is near P and P'. The pressure gradient acts on the fluid in the boundary layer, accelerating it between P and Q but decelerating it between Q and P'. As the flow velocity increases, so does the pressure gradient, and at a certain stage the decelerating effect between Q and P' becomes so large that the direction of flow within the boundary layer reverses in sign near the point labeled R in the diagram. The backflow of fluid near R causes an accumulation of fluid that obliges the oncoming boundary layer to separate, and the fluid behind the sphere circulates slowly within the boundary layer as a ring-shaped eddy (Figure 58B).

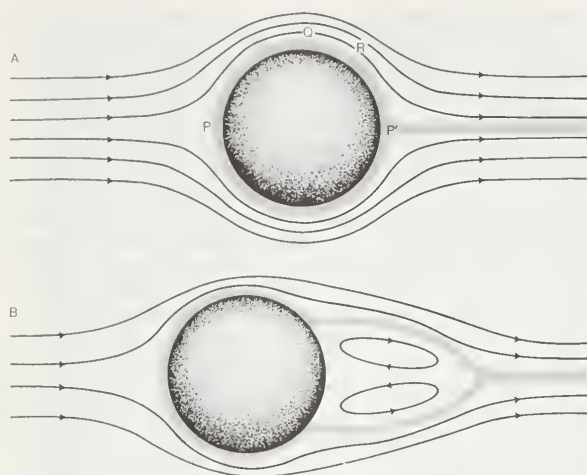


Figure 58: Flow past a stationary solid sphere. The fluid is free of vorticity outside a boundary layer, which is represented by shading. In (A) the boundary layer is still attached to the sphere, though it continues downstream of it. In (B) it has separated, and an eddy has formed behind the sphere.

The diagrams in Figure 58 might well refer to a cylinder rather than a sphere. If such were the case, however, the regions of circulating flow behind the obstacle that are shown in the second diagram would form parts of two separate straight eddies instead of a single ring-shaped one. At high velocities the eddies behind a cylinder become so large that they are blown off by the current and disappear downstream while new eddies form in their place; they are said to have been shed. The top and bottom eddies are shed alternately, and the cylinder experiences an oscillating force as a consequence. If the cylinder is something flexible like a telephone or power cable, it will move to and fro under this force; the singing noise produced by cables in high winds is due to a resonance between their natural frequency of transverse oscillation and the frequency of eddy shedding. Similar processes are liable to occur behind obstacles of any shape, and the occurrence of eddies behind rocks or walls that interrupt the smooth flow of rivers is a familiar phenomenon.

Eddy shedding

DRAG

A fluid stream exerts a drag force F_D on any obstacle placed in its path, and the same force arises if the obstacle moves and the fluid is stationary. How large it is and how it may be reduced are questions of obvious importance to designers of moving vehicles of all sorts and equally to designers of cooling towers and other structures who want to be certain that the structures will not collapse in the face of winds.

An expression for the drag force on a sphere which is valid at such low velocities that the v^2 term in the Navier-Stokes equation is negligible, and thus at velocities such that the boundary layer thickness described by (171) is larger than the sphere diameter D , was first obtained by Stokes. Known as Stokes's law, it may be written as

$$F_D = 3\pi\eta Dv_0 \tag{172}$$

Stokes's law

One-third of this force is transmitted to the sphere by shear stresses near the equator, and the remaining two-thirds are due to the pressure being higher at the front of the sphere than at the rear.

As the velocity increases and the boundary layer decreases in thickness, the effect of the shear stresses (or of what is sometimes called skin friction in this context) becomes less and less important compared with the effect of the pressure difference. It is impossible to calculate that difference precisely, except in the limit to which Stokes's law applies, but there are grounds for expecting that once eddies have formed it is about $\rho v_0^2/2$. Hence at high velocities one may expect

$$F_D = \left(\frac{\rho v_0^2}{2}\right)A', \tag{173}$$

where A' is some effective cross-sectional area, presumably comparable to its true cross-sectional area A (which is $\pi D^2/4$ for a sphere) but not necessarily exactly equal to this. It is conventional to describe drag forces in terms of a dimensionless quantity called the drag coefficient; this is defined, irrespective of the shape of the body, as the ratio $[F_D/(\rho v_0^2/2)A]$ and is denoted by C_D . At high velocities, C_D is clearly the same thing as the ratio (A'/A) and should therefore be of order unity.

This is as far as theory can go with this problem. The principles of dimensional analysis can be invoked to show that, provided the compressibility of the fluid is irrelevant (*i.e.*, provided the flow velocity is well below the speed of sound), the drag coefficient must be some universal function of another dimensionless quantity known as the Reynolds number and defined as

$$R = \frac{\rho v_0 D}{\eta} \tag{174}$$

One must, however, resort to experiments to discover the form of this function. Fortunately, a limited number of experiments will suffice because the function is universal. They can be performed using whatever liquids and spheres are most convenient, provided that the whole range of R that is likely to be important is covered. Once the results have been plotted on a graph of C_D versus R , the graph can be used to predict the drag forces experienced by other spheres in other liquids at velocities that may be quite different from those so far employed. This point is worth emphasizing because it enshrines the principle of dynamic similarity, which is heavily relied on by engineers whenever they use results obtained with models to predict the behaviour of much larger structures.

The C_D versus R curve for spheres, plotted with logarithmic scales, is shown in Figure 59. Stokes's law, re-expressed in terms of C_D and R , becomes $C_D = 24/R$, and it is represented by the straight line on the left of the diagram. This law evidently fails when R exceeds about 1. There is a considerable range of R in the middle of the diagram over which C_D is about 0.5, but when R reaches about 3×10^5 it falls dramatically, to about 0.1. The figure includes the corresponding curves for cylinders of diameter D whose axes are transverse to the direction of flow and for transverse disks of diameter D . The curve for cylinders is similar to that for spheres (though it has no straight-line part at low Reynolds number to correspond to Stokes's law), but the curve for disks is noticeably flatter. This flatness is linked to the fact that a disk has sharp edges around which the streamlines converge and diverge rapidly. The resulting large pressure gradients near the edge favour the formation and shedding of eddies. The drag force on a transverse flat plate of any shape can normally be estimated quite accurately, provided its edges are sharp, by assuming the drag coefficient to be unity.

Since sharp edges favour the formation and shedding of eddies, and thereby increase the drag coefficient, one may hope to reduce the drag coefficient by streamlining the obstacle. It is at the rear of the obstacle that separation occurs, and it is therefore the rear that needs streamlining.

Reducing drag

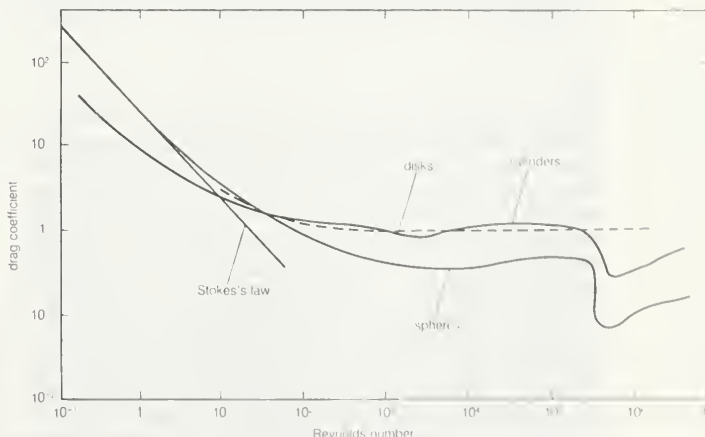


Figure 59: Variation of drag coefficient with Reynolds number for spheres, cylinders, and disks (see text).

By stretching this out in the manner suggested in Figure 60A, the pressure gradient acting on the boundary layer behind the obstacle can be much reduced. Other methods of reducing drag that have some practical applications are illustrated in Figures 60B and 60C. In 60B the obstacle is the wing of an aircraft with a slot through its leading edge; the current of air channeled through this slot imparts forward momentum to the fluid in the boundary layer on the upper surface of the wing to hinder this fluid from moving backward. The cowls that are often fitted to the leading edges of aircraft wings have a similar purpose. In Figure 60C, the obstacle is equipped with an internal device—a pump of some sort—which prevents the accumulation of boundary-layer fluid that would otherwise lead to separation by sucking it in through small holes in the surface of the obstacle, near Q; the fluid may be ejected again through holes near P', where it will do no harm.

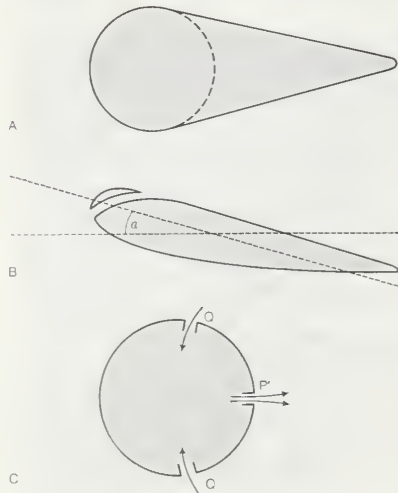


Figure 60: *Methods for reducing drag.* Each diagram represents a solid object that is stationary in the path of a fluid flowing from left to right, or that is moving from right to left through a fluid which is stationary: (A) a sphere that has been streamlined; (B) an aircraft wing, inclined at angle α , which is slotted along its leading edge; (C) a sphere equipped with an internal pump, which draws fluid in near Q and expels it at P'.

It should be stressed that the curves in Figure 59 are universal only so long as the velocity v_0 is much less than the speed of sound. When v_0 is comparable with the speed of sound, V_s , the compressibility of the fluid becomes relevant, which means that the drag coefficient has to be regarded as dependent on the dimensionless ratio $M = v_0/V_s$, known as the Mach number, as well as on the Reynolds number. The drag coefficient always rises as M approaches unity but may thereafter fall. To reduce drag in the supersonic region, it pays to streamline the front of obstacles or projectiles rather than the rear, as this reduces the intensity of the shock cone (see above *Compressible flow in gases*).

LIFT

If an aircraft wing, or airfoil, is to fulfill its function, it must experience an upward lift force, as well as a drag force, when the aircraft is in motion. The lift force arises because the speed at which the displaced air moves over the top of the airfoil (and over the top of the attached boundary layer) is greater than the speed at which it moves over the bottom and because the pressure acting on the airfoil from below is therefore greater than the pressure from above. It also can be seen, however, as an inevitable consequence of the finite circulation that exists around the airfoil. One way to establish circulation around an obstacle is to rotate it, as was seen earlier in the description of the Magnus effect. The circulation around an airfoil, however, is created by its forward motion; it arises as soon as the airfoil moves fast enough to shed its first eddy.

The lift force on an airfoil moving through stationary

air at a steady speed v_0 is the same as the lift force on an identical airfoil that is stationary in air moving at v_0 the other way; the latter is easier to represent pictorially. Figure 61A shows a set of streamlines representing potential flow past a stationary inclined plate before any eddy has been shed. The pattern is a symmetrical one, and the pressure variations associated with it generate neither drag nor lift. At the rear of the plate, however, the streamlines diverge rapidly, so conditions exist for the formation of an eddy there, and the sense of its rotation will be counterclockwise. It grows more easily and is shed more quickly because the edges of the plate are sharp. Figure 61B shows some streamlines for the same plate a moment after shedding when the detached eddy, known as the starting vortex, is still in view. The circulation around the closed loop shown by a broken curve in this diagram was zero before the eddy formed and, according to Thomson's theorem (see above), it must still be zero. Passing through this loop, there thus must be a vortex line having clockwise circulation $-K$ to compensate for the circulation $+K$ of the starting vortex. This other line, known as the bound vortex, is not immediately apparent in the diagram because it is attached to the plate, and it remains thus attached as the starting vortex is swept away downstream. It does show up, however, in a modification of the flow pattern immediately behind the plate, where the streamlines no longer diverge as they do in Figure 61A. Because the divergence here has been eliminated, no further eddies are likely to be formed.

Earlier, the formula $\rho v_0 K$ was quoted for the strength of the Magnus force per unit length of a rotating cylinder, and the same formula can be applied to the inclined plate in Figure 61B or to any airfoil that has shed a starting vortex and around which, consequently, there is circulation. The validity of the formula does not depend in any way on the precise shape of the airfoil, any more than the force exerted by a magnetic field on a wire carrying a current depends on the cross-sectional shape of the wire. The design of the airfoil, nevertheless, has a critical effect on the magnitude of the lift force because it determines the magnitude of K . The sort of cross section that is adopted for the wings of aircraft has been sketched already in Figure 60B. The rear edge is made as sharp as possible for reasons that have already been explained, and it may take the form of hinged flaps that are lowered at takeoff. Lowering the flaps increases K and therefore also the lift, but the flaps need to be raised when the aircraft has reached its cruising altitude because they cause undesirable drag. The circulation and the lift can also be increased by increasing

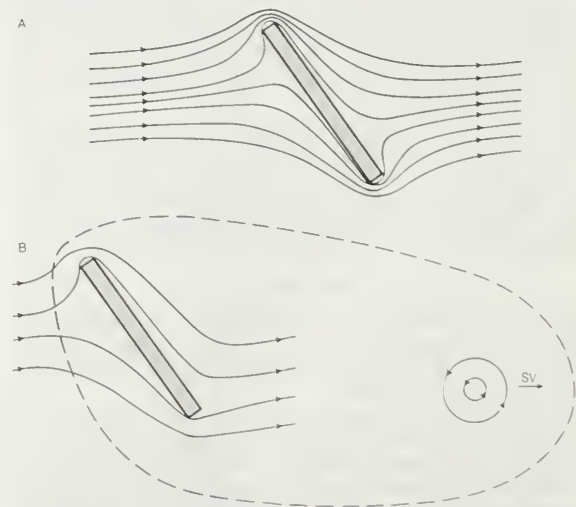


Figure 61: *Generation of lift force.* (A) Streamlines for potential flow past a stationary inclined plate; conditions exist for formation of an eddy behind the plate, with counterclockwise rotation. (B) An eddy formed behind the plate in (A) has been shed as a starting vortex (SV), which is being carried downstream. The circulation remains zero around the large loop indicated by a broken line, and the streamlines no longer diverge behind the rear edge of the plate, as they do in (A) (see text).

Starting vortex and bound vortex

the angle α (see Figure 60B) at which the main part of the airfoil is inclined to the direction of motion. There is a limit to the lift that can be generated in this way, however, for if the inclination is too great the boundary layer separates behind the wing's leading edge, and the bound vortex, on which the lift depends, may be shed as a result. The aircraft is then said to stall. The leading edge is made as smooth and rounded as possible to discourage stalling.

Thomson's theorem can be used to prove that if the airfoil is of finite length then the starting vortex and the bound vortex must both be parts of a single, continuous vortex ring. They are joined by two trailing vortices, which run backward from the ends of the airfoil. As time passes, these trailing vortices grow steadily longer, and more and more energy is needed to feed the swirling motion of the fluid around them. It is clear, at any rate in the case where the airfoil is moving and the air is stationary, that this energy can come only from whatever agency propels the airfoil forward, and hence that the trailing vortices are a source of additional drag. The magnitude of the additional drag is proportional to K^2 but it does not increase, as the lift force does, if the airfoil is made longer while K is kept the same. For this reason, designers who wish to maximize the ratio of lift to drag will make the wings of their aircraft as long as they can—as long, that is, as is consistent with strength and rigidity requirements.

When a yacht is sailing into the wind, its sail acts as an airfoil of which the mast is the leading edge, and the considerations that favour long wings for aircraft favour tall masts as well.

TURBULENCE

The nonlinear nature of the $(\mathbf{v} \cdot \nabla)\mathbf{v}$ term in the Navier-Stokes equation—equation (155)—means that solutions of this equation cannot be superposed. The fact that $\mathbf{v}_1(\mathbf{R}, t)$ and $\mathbf{v}_2(\mathbf{R}, t)$ satisfy the equation does not ensure that $(\mathbf{v}_1 + \mathbf{v}_2)$ does so too. The nonlinear term provides a contact, in fact, through which two different modes of motion may exchange energy, so that one grows in amplitude at the expense of the other. A great deal of experimental and theoretical work has shown, in particular, that if a fluid is undergoing regular laminar motion (of the sort that was discussed in connection with Poiseuille's law, for example) at sufficiently high rates of shear, small periodic perturbations of this motion are liable to grow parasitically. Perturbations on a smaller scale still grow parasitically on those that are first established, until the flow pattern is so grossly disturbed that it is no longer useful to define a fluid velocity for each point in space; the description of the flow has to be a statistical one in terms of mean values and of correlated fluctuations about the mean. The flow is then said to be turbulent.

In the case (to which Poiseuille's law applies) of laminar flow through a uniform cylindrical pipe of diameter D , turbulence inevitably sets in when the Reynolds number R reaches a critical value that is about 10^3 ; in this context, the Reynolds number is defined (compare equation [174]) as

$$R = \frac{4\rho Q}{\pi D\eta} = \frac{\rho D \langle v \rangle}{\eta}, \tag{175}$$

where Q is the rate of discharge and $\langle v \rangle$ is the mean fluid velocity. Turbulence sets in at much lower velocities, however, if the end of the pipe where the fluid enters is not carefully flared. The critical value of the Reynolds number for a pipe with a bluff entry may be as low as 2300, and this corresponds to a rate of discharge through a pipe for which D is, say, two centimetres, of only about three litres per minute. Thus pipe flow in engineering practice is more often turbulent than not. Once turbulence has set in, Q increases less rapidly with pressure gradient than Poiseuille's equation—equation (150)—predicts; it increases roughly as the square root of the pressure gradient or slightly more rapidly than this if the internal surface of the pipe is very smooth.

Turbulence arises not only in pipes but also within boundary layers around solid obstacles when the rate of shear within the boundary layer becomes large enough. Curiously enough, the onset of turbulence in the boundary

layer can reduce the drag force on obstacles. In the case of a spherical obstacle, the point at which the boundary layer separates from the rear surface of the sphere shifts backward when the boundary layer becomes turbulent, away from the equator Q in Figure 58 and toward P' , and the eddies attached to the sphere therefore become smaller. It is turbulence in the boundary layer that is responsible for the dramatic drop in the drag coefficient for both spheres and cylinders that occurs, as can be seen from Figure 59, when the Reynolds number is about 3×10^5 . This drop enables golf balls to travel farther than they would do otherwise, and the dimples on the surface of golf balls are meant to encourage turbulence in the boundary layer. If swimsuits with rough surfaces help swimmers to move faster, as has been claimed, the same explanation may apply.

Where conditions for turbulence exist, flow rates of water through tubes may be increased and the drag forces exerted on obstacles by water diminished by dissolving small amounts of suitable polymers in the water. This is surprising, because such additives increase viscosity, and in the preturbulent regime to which Poiseuille's law applies, their effect on the flow rate is quite the reverse. As has already been stated, the small perturbations that arise in a turbulent fluid tend to collapse into smaller perturbations and then into smaller perturbations still, until the motion is turbulent on a very fine scale—i.e., on the scale of molecular dimensions—and until the energy stored in the perturbations is finally dissipated as heat. Polymer molecules seem to have the effect they do because, over the relatively large distances to which each such molecule extends, they impose a coherence on the fluid motion that would not otherwise be present.

CONVECTION

Apart from some remarks in the above section *Compressible flow in gases* about the circulation of the atmosphere, no attention has yet been paid to situations in which temperature differences are imposed upon a fluid by contact with hot and cold bodies. This subject will be briefly taken up here.

Consider first the case of two vertical plates with fluid between them, one at temperature T_1 and the other at T_2 , in the presence of a vertical gravitational field. The hotter plate might be a domestic radiator and the colder plate the wall to which it is fixed. Thermal conduction ensures that the layer of air adjacent to the radiator is hotter than the rest of the air, and thermal expansion ensures that it is less dense. Consequently, the vertical pressure gradient which satisfies equation (123) in the rest of the air is too large to keep the layer adjacent to the radiator in equilibrium; that layer rises and, similarly, the cold layer adjacent to the wall falls. A circulating pattern of thermal convection is thereby established, and, because this brings colder air into contact with the radiator, the rate at which heat is lost from the radiator is enhanced. The heat loss, once convection has been established, depends in a complicated manner on the separation between the plates (D) and on the thermal diffusivity (κ), specific heat, density, thermal expansion coefficient (α), and viscosity of the fluid. The heat loss also depends on $(T_1 - T_2)$, of course, and it is worthwhile noting that the manner in which it does so is not linear; the heat loss increases more rapidly than the temperature difference. Newton's law of cooling, which postulates a linear relationship, is obeyed only in circumstances where convection is prevented or in circumstances where it is forced (when a radiator is fan-assisted, for example).

Imagine a situation in which the same two plates are horizontal rather than vertical. In such a case, no convection can occur if the hot plate is above the cold one, and it is not obvious that it occurs in the reverse situation. Whether it does so or not depends on the magnitude of the temperature difference through a dimensionless combination of some of the relevant parameters, $\rho g \alpha D^3 (T_1 - T_2) / \eta \kappa$, which is known as the Rayleigh number. If the Rayleigh number is less than 1,708, the fluid is stable—or perhaps it would be more accurate to say that it is metastable—even though it is warmer at the bottom than at the top. How-

Trailing vortices

Turbulent flow

Turbulence within boundary layers

Effect of polymer additives

Bénard
cells

ever, when 1,708 is exceeded, a pattern of convective rolls known as Bénard cells is established between the plates. Evidence for the existence of such cells in the convecting atmosphere is sometimes seen in the regular columns of cloud that form over regions where the air is rising. Their periodicity can be astonishingly uniform.

Macroscopic instabilities of a convective nature, of which

the formation of Bénard cells provides just one example, are a feature of the oceans as well as of the atmosphere and are frequently associated with gradients of salinity rather than gradients of temperature. A serious discussion of atmospheric and oceanic circulation on the Earth, however, requires a more detailed examination of the dynamics of rotating fluids than is given here. (T.E.F.)

QUANTUM MECHANICS

Quantum mechanics deals with the behaviour of matter and light on the atomic and subatomic scale. It attempts to describe and account for the properties of molecules and atoms and their constituents—electrons, protons, neutrons, and other more esoteric particles such as quarks and gluons. These properties include the interactions of the particles with one another and with electromagnetic radiation (*i.e.*, light, X rays, and gamma rays).

The behaviour of matter and radiation on the atomic scale often seems peculiar, and the consequences of quantum theory are accordingly difficult to understand and to believe. Its concepts frequently conflict with commonsense notions derived from observations of the everyday world. There is no reason, however, why the behaviour of the atomic world should conform to that of the familiar, large-scale world. It is important to realize that quantum mechanics is a branch of physics and that the business of physics is to describe and account for the way the world—on both the large and the small scale—actually is and not how one imagines it or would like it to be.

The study of quantum mechanics is rewarding for several reasons. First, it illustrates the essential methodology of physics. Second, it has been enormously successful in giving correct results in practically every situation to which it has been applied. There is, however, an intriguing paradox. In spite of the overwhelming practical success of quantum mechanics, the foundations of the subject contain unresolved problems—in particular, problems concerning the nature of measurement. An essential feature of quantum mechanics is that it is generally impossible, even in principle, to measure a system without disturbing it; the detailed nature of this disturbance and the exact point at which it occurs are obscure and controversial. Thus, quantum mechanics has attracted some of the ablest scientists of the 20th century, and they have erected what is perhaps the finest intellectual edifice of the period.

Historical basis of quantum theory

BASIC CONSIDERATIONS

At a fundamental level, both radiation and matter have characteristics of particles and waves. The gradual recognition by scientists that radiation has particle-like properties and that matter has wavelike properties provided the impetus for the development of quantum mechanics. Influenced by Newton, most physicists of the 18th century believed that light consisted of particles, which they called corpuscles. From about 1800, evidence began to accumulate for a wave theory of light. At about this time Thomas Young showed that, if monochromatic light passes through a pair of slits, the two emerging beams interfere, so that a fringe pattern of alternately bright and dark bands appears on a screen. The bands are readily explained by a wave theory of light. According to the theory, a bright band is produced when the crests (or troughs) of the waves from the two slits arrive together at the screen; a dark band is produced when the crest of one wave arrives at the same time as the trough of the other, and the effects of the two light beams cancel. Beginning in 1815, a series of experiments by Augustin-Jean Fresnel of France and others showed that, when a parallel beam of light passes through a single slit, the emerging beam is no longer parallel but starts to diverge; this phenomenon is known as diffraction. Given the wavelength of the light and the geometry of the apparatus (*i.e.*, the separation and widths of the slits and the distance from the slits to the screen), one can use the wave theory to calculate the

expected pattern in each case; the theory agrees precisely with the experimental data.

EARLY DEVELOPMENTS

Planck's radiation law. By the end of the 19th century, physicists almost universally accepted the wave theory of light. However, though the ideas of classical physics explain interference and diffraction phenomena relating to the propagation of light, they do not account for the absorption and emission of light. All bodies radiate electromagnetic energy as heat; in fact, a body emits radiation at all wavelengths. The energy radiated at different wavelengths is a maximum at a wavelength that depends on the temperature of the body; the hotter the body, the shorter the wavelength for maximum radiation. Attempts to calculate the energy distribution for the radiation from a blackbody using classical ideas were unsuccessful. (A blackbody is a hypothetical ideal body or surface that absorbs and reemits all radiant energy falling on it.) One formula, proposed by Wilhelm Wien of Germany, did not agree with observations at long wavelengths, and another, proposed by Lord Rayleigh (John William Strutt) of England, disagreed with those at short wavelengths.

In 1900 the German theoretical physicist Max Planck made a bold suggestion. He assumed that the radiation energy is emitted, not continuously, but rather in discrete packets called quanta. The energy E of the quantum is related to the frequency ν by

$$E = h\nu. \quad (176)$$

The quantity h , now known as the Planck constant, is a universal constant with the approximate value $h = 6.63 \times 10^{-34}$ joule-second. (The more precise, currently accepted value is 6.626075×10^{-34} joule-second.) Planck showed that the calculated energy spectrum then agreed with observation over the entire wavelength range.

Einstein and the photoelectric effect. In 1905 Einstein extended Planck's hypothesis to explain the photoelectric effect, which is the emission of electrons by a metal surface when it is irradiated by light or X rays. The kinetic energy of the emitted electrons depends on the frequency ν of the radiation, not on its intensity; for a given metal, there is a threshold frequency ν_0 below which no electrons are emitted. Furthermore, emission takes place as soon as the light shines on the surface; there is no detectable delay. Einstein showed that these results can be explained by two assumptions: (1) that light is composed of corpuscles or photons, the energy of which is given by Planck's relationship, and (2) that an atom in the metal can absorb either a whole photon or nothing. Part of the energy of the absorbed photon frees an electron, which requires a fixed energy W , known as the work function of the metal; the rest is converted into the kinetic energy $\frac{1}{2}m_e u^2$ of the emitted electron (m_e is the mass of the electron and u is its velocity). Thus, the energy relation is

$$h\nu = W + \frac{1}{2}m_e u^2. \quad (177)$$

If ν is less than ν_0 , where $h\nu_0 = W$, no electrons are emitted. Not all the experimental results mentioned above were known in 1905, but all Einstein's predictions have been verified since.

Bohr's theory of the atom. A major contribution to the subject was made by Niels Bohr of Denmark, who applied the quantum hypothesis to atomic spectra in 1913. The spectra of light emitted by gaseous atoms had been studied extensively since the mid-19th century. It was found that

Concept of
quantumSupport
for a wave
theory of
light

radiation from gaseous atoms at low pressure consists of a set of discrete wavelengths. This is quite unlike the radiation from a solid, which is distributed over a continuous range of wavelengths. The set of discrete wavelengths from gaseous atoms is known as a line spectrum, because the image of the linear slit in the spectrometer is a series of sharp lines. The wavelengths of the lines are characteristic of the element and may form extremely complex patterns. The simplest spectra are those of atomic hydrogen and the alkali atoms (*e.g.*, lithium, sodium, and potassium). For hydrogen, the wavelengths λ are given by the empirical formula

$$\frac{1}{\lambda} = R_{\infty} \left(\frac{1}{m^2} - \frac{1}{n^2} \right), \quad (178)$$

where m and n are positive integers with $n > m$ and R_{∞} , known as the Rydberg constant, has the value 1.0973731×10^7 per metre. For a given value of m , the lines for varying n form a series. The lines for $m = 1$, the Lyman series, lie in the ultraviolet part of the spectrum. Those for $m = 2$, the Balmer series, lie in the visible spectrum, and so on.

Bohr started with a model suggested by the New Zealand-born British physicist Ernest Rutherford. The model was based on the experiments of Hans Geiger and Ernest Marsden, who in 1909 bombarded gold atoms with massive, fast-moving alpha particles; when some of these particles were deflected backward, Rutherford concluded that the atom has a massive, charged nucleus. In Rutherford's model, the atom resembles a miniature solar system with the nucleus acting as the Sun and the electrons as the circulating planets. Bohr made three assumptions. First, he postulated that, in contrast to classical mechanics, where an infinite number of orbits is possible, an electron can be in only one of a discrete set of orbits, which he termed stationary states. Second, he postulated that the only orbits allowed are those for which the angular momentum of the electron is a whole number n times \hbar (\hbar stands for $h/2\pi$). Third, Bohr assumed that Newton's laws of motion, so successful in calculating the paths of the planets around the Sun, also applied to electrons orbiting the nucleus. The force on the electron (the analogue of the gravitational force between the Sun and a planet) is the electrostatic attraction between the positively charged nucleus and the negatively charged electron. With these simple assumptions, he showed that the energy of the orbit has the form

$$E_n = -\frac{E_0}{n^2}, \quad (179)$$

where E_0 is a constant that may be expressed by a combination of the known constants e , m_e , and \hbar . While in a stationary state, the atom does not give off energy as light; however, when an electron makes a transition from a state with energy E_n to one with lower energy E_m , a quantum of energy is radiated with frequency ν , given by the equation

$$h\nu = E_n - E_m. \quad (180)$$

Inserting the expression for E_n into this equation and using the relation $\lambda\nu = c$, where c is the speed of light, Bohr derived the formula for the wavelengths of the lines in the hydrogen spectrum, with the correct value of the Rydberg constant.

Bohr's theory was a brilliant step forward. Its two most important features have survived in present-day quantum mechanics. They are (1) the existence of stationary, nonradiating states and (2) the relationship of radiation frequency to the energy difference between the initial and final states in a transition. Prior to Bohr, physicists had thought that the radiation frequency would be the same as the electron's frequency of rotation in an orbit.

Scattering of X rays. Soon scientists were faced with the fact that another form of radiation, X rays, also exhibits both wave and particle properties. Max von Laue of Germany had shown in 1912 that crystals can be used as three-dimensional diffraction gratings for X rays; his technique constituted the fundamental evidence for the wavelike nature of X rays. The atoms of a crystal, which

are arranged in a regular lattice, scatter the X rays. For certain directions of scattering, all the crests of the X rays coincide. (The scattered X rays are said to be in phase and to give constructive interference.) For these directions, the scattered X-ray beam is very intense. Clearly, this phenomenon demonstrates wave behaviour. In fact, given the interatomic distances in the crystal and the directions of constructive interference, the wavelength of the waves can be calculated.

In 1922 the American physicist Arthur Holly Compton showed that X rays scatter from electrons as if they are particles. Compton performed a series of experiments on the scattering of monochromatic, high-energy X rays by graphite. He found that part of the scattered radiation had the same wavelength λ_0 as the incident X rays but that there was an additional component with a longer wavelength λ . To interpret his results, Compton regarded the X-ray photon as a particle that collides and bounces off an electron in the graphite target as though the photon and the electron were a pair of (dissimilar) billiard balls. Application of the laws of conservation of energy and momentum to the collision leads to a specific relation between the amount of energy transferred to the electron and the angle of scattering. For X rays scattered through an angle θ , the wavelengths λ and λ_0 are related by the equation

$$\lambda - \lambda_0 = \frac{h}{m_e c} (1 - \cos \theta). \quad (181)$$

The experimental correctness of Compton's formula is direct evidence for the corpuscular behaviour of radiation.

Brogie's wave hypothesis. Faced with evidence that electromagnetic radiation has both particle and wave characteristics, Louis-Victor de Broglie of France suggested a great unifying hypothesis in 1924. Broglie proposed that matter has wave, as well as particle, properties. He suggested that material particles can behave as waves and that their wavelength λ is related to the linear momentum p of the particle by

$$\lambda = \frac{h}{p}. \quad (182)$$

In 1927 Clinton Davisson and Lester Germer of the United States confirmed Broglie's hypothesis for electrons. Using a crystal of nickel, they diffracted a beam of monoenergetic electrons and showed that the wavelength of the waves is related to the momentum of the electrons by the Broglie equation. Since Davisson and Germer's investigation, similar experiments have been performed with atoms, molecules, neutrons, protons, and many other particles. All behave like waves with the same wavelength-momentum relationship.

Basic concepts and methods

Bohr's theory, which assumed that electrons moved in circular orbits, was extended by the German physicist Arnold Sommerfeld and others to include elliptic orbits and other refinements. Attempts were made to apply the theory to more complicated systems than the hydrogen atom. However, the ad hoc mixture of classical and quantum ideas made the theory and calculations increasingly unsatisfactory. Then in the 12 months starting in July 1925, a period of creativity without parallel in the history of physics, there appeared a series of papers by German scientists that set the subject on a firm conceptual foundation. The papers took two approaches: (1) matrix mechanics, proposed by Werner Heisenberg, Max Born, and Pascual Jordan, and (2) wave mechanics, put forward by Erwin Schrödinger. The protagonists were not always polite to each other. Heisenberg found the physical ideas of Schrödinger's theory "disgusting," and Schrödinger was "discouraged and repelled" by the lack of visualization in Heisenberg's method. However, Schrödinger, not allowing his emotions to interfere with his scientific endeavours, showed that, in spite of apparent dissimilarities, the two theories are equivalent mathematically. The present discussion follows Schrödinger's wave mechanics because it is less abstract and easier to understand than Heisenberg's matrix mechanics.

Line spectrum

Most important features of Bohr's theory

Wave-particle hypothesis

SCHRÖDINGER'S WAVE MECHANICS

Schrödinger expressed Broglie's hypothesis concerning the wave behaviour of matter in a mathematical form that is adaptable to a variety of physical problems without additional arbitrary assumptions. He was guided by a mathematical formulation of optics, in which the straight-line propagation of light rays can be derived from wave motion when the wavelength is small compared to the dimensions of the apparatus employed. In the same way, Schrödinger set out to find a wave equation for matter that would give particle-like propagation when the wavelength becomes comparatively small. According to classical mechanics, if a particle of mass m_e is subjected to a force such that its potential energy is $V(x,y,z)$ at position x,y,z , then the sum of $V(x,y,z)$ and the kinetic energy $p^2/2m_e$ is equal to a constant, the total energy E of the particle. Thus,

$$\frac{p^2}{2m_e} + V(x,y,z) = E. \quad (183)$$

It is assumed that the particle is bound—*i.e.*, confined by the potential to a certain region in space because its energy E is insufficient for it to escape. Since the potential varies with position, two other quantities do also: the momentum and, hence, by extension from the Broglie relation, the wavelength of the wave. Postulating a wave function $\Psi(x,y,z)$ that varies with position, Schrödinger replaced p in the above energy equation with a differential operator that embodied the Broglie relation. He then showed that Ψ satisfies the partial differential equation

$$-\frac{\hbar^2}{2m_e} \left(\frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} + \frac{\partial^2 \Psi}{\partial z^2} \right) + V(x,y,z)\Psi = E\Psi. \quad (184)$$

Time-independent Schrödinger wave equation

This is the (time-independent) Schrödinger wave equation of 1926, which established quantum mechanics in a widely applicable form. An important advantage of Schrödinger's theory is that no further arbitrary quantum conditions need be postulated. The required quantum results follow from certain reasonable restrictions placed on the wave function—for example, that it should not become infinitely large at large distances from the centre of the potential.

Schrödinger applied his equation to the hydrogen atom, for which the potential function, given by classical electrostatics, is proportional to $-e^2/r$, where $-e$ is the charge on the electron. The nucleus (a proton of charge e) is situated at the origin, and r is the distance from the origin to the position of the electron. Schrödinger solved the equation for this particular potential with straightforward, though not elementary, mathematics. Only certain discrete values of E lead to acceptable functions Ψ . These functions are characterized by a trio of integers n, l, m , termed quantum numbers. The values of E depend only on the integers n (1, 2, 3, etc.) and are identical with those given by the Bohr theory. The quantum numbers l and m are related to the angular momentum of the electron; $\sqrt{l(l+1)}\hbar$ is the magnitude of the angular momentum, and $m\hbar$ is its component along some physical direction.

The square of the wave function, Ψ^2 , has a physical interpretation. Schrödinger originally supposed that the electron was spread out in space and that its density at point x,y,z was given by the value of Ψ^2 at that point. Almost immediately Born proposed what is now the accepted interpretation—namely, that Ψ^2 gives the probability of finding the electron at x,y,z . The distinction between the two interpretations is important. If Ψ^2 is small at a particular position, the original interpretation implies that a small fraction of an electron will always be detected there. In Born's interpretation, nothing will be detected there most of the time, but, when something is observed, it will be a whole electron. Thus, the concept of the electron as a point particle moving in a well-defined path around the nucleus is replaced in wave mechanics by clouds that describe the probable locations of electrons in different states.

ELECTRON SPIN AND ANTIPARTICLES

In 1928 the English physicist Paul A.M. Dirac produced a wave equation for the electron that combined relativity with quantum mechanics. Schrödinger's wave equation

does not satisfy the requirements of the special theory of relativity because it is based on a nonrelativistic expression for the kinetic energy ($p^2/2m_e$). Dirac showed that an electron has an additional quantum number m_s . Unlike the first three quantum numbers, m_s is not a whole integer and can have only the values $+1/2$ and $-1/2$. It corresponds to an additional form of angular momentum ascribed to a spinning motion. (The angular momentum mentioned above is due to the orbital motion of the electron, not its spin.) The concept of spin angular momentum was introduced in 1925 by Samuel A. Goudsmit and George E. Uhlenbeck, two graduate students at the University of Leiden, Neth., to explain the magnetic moment measurements made by Otto Stern and Walther Gerlach of Germany several years earlier. The magnetic moment of a particle is closely related to its angular momentum; if the angular momentum is zero, so is the magnetic moment. Yet Stern and Gerlach had observed a magnetic moment for electrons in silver atoms, which were known to have zero orbital angular momentum. Goudsmit and Uhlenbeck proposed that the observed magnetic moment was attributable to spin angular momentum.

The electron-spin hypothesis not only provided an explanation for the observed magnetic moment but also accounted for many other effects in atomic spectroscopy, including changes in spectral lines in the presence of a magnetic field (Zeeman effect), doublet lines in alkali spectra, and fine structure (close doublets and triplets) in the hydrogen spectrum.

The Dirac equation also predicted additional states of the electron that had not yet been observed. Experimental confirmation was provided in 1932 by the discovery of the positron by the American physicist Carl David Anderson. Every particle described by the Dirac equation has to have a corresponding antiparticle, which differs only in charge. The positron is just such an antiparticle of the negatively charged electron, having the same mass as the latter but a positive charge.

IDENTICAL PARTICLES AND MULTIELECTRON ATOMS

Because electrons are identical to (*i.e.*, indistinguishable from) each other, the wave function of an atom with more than one electron must satisfy special conditions. The problem of identical particles does not arise in classical physics, where the objects are large-scale and can always be distinguished, at least in principle. There is no way, however, to differentiate two electrons in the same atom, and the form of the wave function must reflect this fact. The overall wave function Ψ of a system of identical particles depends on the coordinates of all the particles. If the coordinates of two of the particles are interchanged, the wave function must remain unaltered or, at most, undergo a change of sign; the change of sign is permitted because it is Ψ^2 that occurs in the physical interpretation of the wave function. If the sign of Ψ remains unchanged, the wave function is said to be symmetric with respect to interchange; if the sign changes, the function is antisymmetric.

The symmetry of the wave function for identical particles is closely related to the spin of the particles. In quantum field theory (see below *Quantum electrodynamics*), it can be shown that particles with half-integral spin ($1/2, 3/2$, etc.) have antisymmetric wave functions. They are called fermions after the Italian-born physicist Enrico Fermi. Examples of fermions are electrons, protons, and neutrons, all of which have spin $1/2$. Particles with zero or integral spin (*e.g.*, mesons, photons) have symmetric wave functions and are called bosons after the Indian mathematician and physicist Satyendra Nath Bose, who first applied the ideas of symmetry to photons in 1924–25.

The requirement of antisymmetric wave functions for fermions leads to a fundamental result, known as the exclusion principle, first proposed in 1925 by the Austrian physicist Wolfgang Pauli. The exclusion principle states that two fermions in the same system cannot be in the same quantum state. If they were, interchanging the two sets of coordinates would not change the wave function at all, which contradicts the result that the wave function must change sign. Thus, two electrons in the same atom cannot have an identical set of values for the four quan-

Dirac's relativistic wave equation

Discovery of the positron

Exclusion principle

tum numbers n, l, m, m_s . The exclusion principle forms the basis of many properties of matter, including the periodic classification of the elements, the nature of chemical bonds, and the behaviour of electrons in solids; the last determines in turn whether a solid is a metal, insulator, or semiconductor (see ATOMS; MATTER).

The Schrödinger equation cannot be solved precisely for atoms with more than one electron. The principles of the calculation are well understood, but the problems are complicated by the number of particles and the variety of forces involved. The forces include the electrostatic forces between the nucleus and the electrons and between the electrons themselves, as well as weaker magnetic forces arising from the spin and orbital motions of the electrons. Despite these difficulties, approximation methods introduced by the English physicist Douglas R. Hartree and others in the 1920s have achieved considerable success. Such schemes start by assuming that each electron moves independently in an average electric field because of the nucleus and the other electrons—*i.e.*, correlations between the positions of the electrons are ignored. Each electron has its own wave function, called an orbital. The overall wave function for all the electrons in the atom satisfies the exclusion principle. Corrections to the calculated energies are then made, which depend on the strengths of the electron-electron correlations and the magnetic forces.

Hartree method

TIME-DEPENDENT SCHRÖDINGER EQUATION

At the same time that Schrödinger proposed his time-independent equation to describe the stationary states, he also proposed a time-dependent equation to describe how a system changes from one state to another. By replacing the energy E in equation (184) with a time-derivative operator, he generalized his wave equation to determine the time variation of the wave function as well as its spatial variation. The time-dependent Schrödinger equation reads

$$-\frac{\hbar^2}{2m_e} \left(\frac{\partial^2 \Psi}{\partial x^2} + \frac{\partial^2 \Psi}{\partial y^2} + \frac{\partial^2 \Psi}{\partial z^2} \right) + V(x, y, z) \Psi = i\hbar \frac{\partial \Psi}{\partial t} \quad (185)$$

The quantity i is the square root of -1 . The function Ψ varies with time t as well as with position x, y, z . For a system with constant energy, E , Ψ has the form

$$\Psi(x, y, z, t) = \Psi(x, y, z) \exp\left(-\frac{iEt}{\hbar}\right) \quad (186)$$

where \exp stands for the exponential function, and the time-dependent Schrödinger equation reduces to the time-independent form.

The probability of a transition between one atomic stationary state and some other state can be calculated with the aid of the time-dependent equation. For example, an atom may change spontaneously from one state to another state with less energy, emitting the difference in energy as a photon with a frequency given by the Bohr relation. If electromagnetic radiation is applied to a set of atoms and if the frequency of the radiation matches the energy difference between two stationary states, transitions can be stimulated. In a stimulated transition, the energy of the atom may increase—*i.e.*, the atom may absorb a photon from the radiation—or the energy of the atom may decrease, with the emission of a photon, which adds to the energy of the radiation. Such stimulated emission processes form the basic mechanism for the operation of lasers. The probability of a transition from one state to another depends on the values of the l, m, m_s quantum numbers of the initial and final states. For most values, the transition probability is effectively zero. However, for certain changes in the quantum numbers, summarized as selection rules, there is a finite probability. For example, according to one important selection rule, the l value changes by unity. The selection rules for radiation relate to the angular momentum properties of the stationary states. The absorbed or emitted photon has its own angular momentum, and the selection rules reflect the conservation of angular momentum between the atoms and the radiation.

Transitions between atomic states

TUNNELING

The phenomenon of tunneling, which has no counterpart in classical physics, is an important consequence of quan-

No counterpart in classical physics

tum mechanics. Consider a particle with energy E in the inner region of a one-dimensional potential well $V(x)$, as shown in Figure 62. (A potential well is a potential that has a lower value in a certain region of space than in the neighbouring regions.) In classical mechanics, if $E < V_0$ (the maximum height of the potential barrier), the particle remains in the well forever; if $E > V_0$, the particle escapes. In quantum mechanics, the situation is not so simple. The particle can escape even if its energy E is below the height of the barrier V_0 , although the probability of escape is small unless E is close to V_0 . In that case, the particle may tunnel through the potential barrier and emerge with the same energy E .

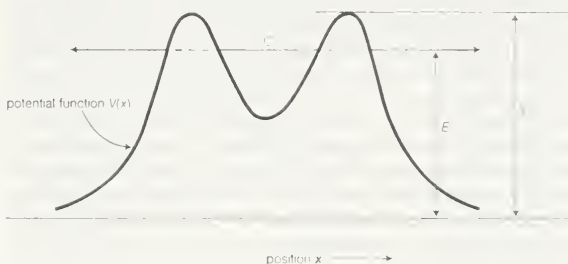


Figure 62: The phenomenon of tunneling. Classically, a particle is bound in the central region C if its energy E is less than V_0 , but in quantum theory the particle may tunnel through the potential barrier and escape.

The phenomenon of tunneling has many important applications. For example, it describes a type of radioactive decay in which a nucleus emits an alpha particle (a helium nucleus). According to the quantum explanation given independently by George Gamow and by Ronald W. Gurney and Edward Condon in 1928, the alpha particle is confined before the decay by a potential of the shape shown in Figure 62. For a given nuclear species, it is possible to measure the energy E of the emitted alpha particle and the average lifetime τ of the nucleus before decay. The lifetime of the nucleus is a measure of the probability of tunneling through the barrier—the shorter the lifetime, the higher the probability. With plausible assumptions about the general form of the potential function, it is possible to calculate a relationship between τ and E that is applicable to all alpha emitters. This theory, which is borne out by experiment, shows that the probability of tunneling, and hence the value of τ , is extremely sensitive to the value of E . For all known alpha-particle emitters, the value of E varies from about 2 to 8 megaelectron volts, or MeV (1 MeV = 10^6 electron volts). Thus, the value of E varies only by a factor of 4, whereas the range of τ is from about 10^{11} years down to about 10^{-6} second, a factor of 10^{17} . It would be difficult to account for this sensitivity of τ to the value of E by any theory other than quantum mechanical tunneling.

AXIOMATIC APPROACH

Although the two Schrödinger equations form an important part of quantum mechanics, it is possible to present the subject in a more general way. Dirac gave an elegant exposition of an axiomatic approach based on observables and states in a classic textbook entitled *The Principles of Quantum Mechanics*. (The book, published in 1930, is still in print.) An observable is anything that can be measured—energy, position, a component of angular momentum, and so forth. Every observable has a set of states, each state being represented by an algebraic function. With each state is associated a number that gives the result of a measurement of the observable. Consider an observable with N states, denoted by $\phi_1, \phi_2, \dots, \phi_N$, and corresponding measurement values a_1, a_2, \dots, a_N . A physical system—*e.g.*, an atom in a particular state—is represented by a wave function Ψ , which can be expressed as a linear combination, or mixture, of the states of the observable. Thus, the Ψ may be written as

$$\Psi = c_1 \phi_1 + c_2 \phi_2 + \dots + c_N \phi_N \quad (187)$$

For a given Ψ , the quantities c_1, c_2 , etc., are a set of numbers that can be calculated. In general, the numbers are

Axiomatic approach based on observables and states

complex, but, in the present discussion, they are assumed to be real numbers.

The theory postulates, first, that the result of a measurement must be an a -value—*i.e.*, a_1 , a_2 , or a_3 , etc. No other value is possible. Second, before the measurement is made, the probability of obtaining the value a_1 is c_1^2 , and that of obtaining the value a_2 is c_2^2 , and so on. If the value obtained is, say, a_5 , the theory asserts that after the measurement the state of the system is no longer the original Ψ but has changed to φ_5 , the state corresponding to a_5 .

A number of consequences follow from these assertions. First, the result of a measurement cannot be predicted with certainty. Only the probability of a particular result can be predicted, even though the initial state (represented by the function Ψ) is known exactly. Second, identical measurements made on a large number of identical systems, all in the identical state Ψ , will produce different values for the measurements. This is, of course, quite contrary to classical physics and common sense, which say that the same measurement on the same object in the same state must produce the same result. Moreover, according to the theory, not only does the act of measurement change the state of the system, but it does so in an indeterminate way. Sometimes it changes the state to φ_1 , sometimes to φ_2 , and so forth.

There is an important exception to the above statements. Suppose that, before the measurement is made, the state Ψ happens to be one of the φ s, say $\Psi = \varphi_3$. Then $c_3 = 1$ and all the other c 's are zero. This means that, before the measurement is made, the probability of obtaining the value a_3 is unity and the probability of obtaining any other value of a is zero. In other words, in this particular case, the result of the measurement can be predicted with certainty. Moreover, after the measurement is made, the state will be φ_3 , the same as it was before. Thus, in this particular case, measurement does not disturb the system. Whatever the initial state of the system, two measurements made in rapid succession (so that the change in the wave function given by the time-dependent Schrödinger equation is negligible) produce the same result.

The value of one observable can be determined by a single measurement. The value of two observables for a given system may be known at the same time, provided that the two observables have the same set of state functions $\varphi_1, \varphi_2, \dots, \varphi_N$. In this case, measuring the first observable results in a state function that is one of the φ s. Because this is also a state function of the second observable, the result of measuring the latter can be predicted with certainty. Thus the values of both observables are known. (Although the φ s are the same for the two observables, the two sets of a values are, in general, different.) The two observables can be measured repeatedly in any sequence. After the first measurement, none of the measurements disturbs the system, and a unique pair of values for the two observables is obtained.

INCOMPATIBLE OBSERVABLES

The measurement of two observables with different sets of state functions is a quite different situation. Measurement of one observable gives a certain result. The state function after the measurement is, as always, one of the states of that observable; however, it is not a state function for the second observable. Measuring the second observable disturbs the system, and the state of the system is no longer one of the states of the first observable. In general, measuring the first observable again does not produce the same result as the first time. To sum up, both quantities cannot be known at the same time, and the two observables are said to be incompatible.

A specific example of this behaviour is the measurement of the component of angular momentum along two mutually perpendicular directions. The Stern-Gerlach experiment mentioned above involved measuring the angular momentum of a silver atom in the ground state. In reconstructing this experiment, a beam of silver atoms is passed between the poles of a magnet. The poles are shaped so that the magnetic field varies greatly in strength over a very small distance (Figure 63). The apparatus determines the m_x quantum number, which can be $+1/2$

or $-1/2$. No other values are obtained. Thus in this case the observable has only two states—*i.e.*, $N = 2$. The inhomogeneous magnetic field produces a force on the silver atoms in a direction that depends on the spin state of the atoms. The result is shown schematically in Figure 64. A beam of silver atoms is passed through magnet A. The atoms in the state with $m_x = +1/2$ are deflected upward and emerge as beam 1, while those with $m_x = -1/2$ are deflected downward and emerge as beam 2. If the direction of the magnetic field is the x -axis, the apparatus measures S_x , which is the x -component of spin angular momentum. The atoms in beam 1 have $S_x = +1/2\hbar$, while those in beam 2 have $S_x = -1/2\hbar$. In a classical picture, these two states represent atoms spinning about the direction of the x -axis with opposite senses of rotation.

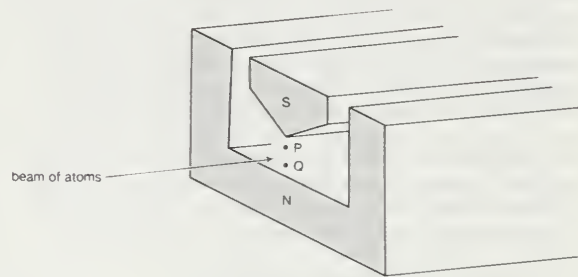


Figure 63: Magnet in Stern-Gerlach experiment. N and S are the north and south poles of a magnet. The knife-edge of S results in a much stronger magnetic field at the point P than at Q.

The y -component of spin angular momentum S_y also can have only the values $+1/2\hbar$ and $-1/2\hbar$; however, the two states of S_y are not the same as for S_x . In fact, each of the states of S_x is an equal mixture of the states for S_y , and conversely, the two S_y states may be pictured as representing atoms with opposite senses of rotation about the y -axis. These classical pictures of quantum states are helpful, but only up to a certain point. For example, quantum theory says that each of the states corresponding to spin about the x -axis is a superposition of the two states with spin about the y -axis. There is no way to visualize this; it has absolutely no classical counterpart. One simply has to accept the result as a consequence of the axioms of the theory. Suppose that, as in Figure 64, the atoms in beam 1 are passed into a second magnet B, which has a magnetic field along the y -axis perpendicular to x . The atoms emerge from B and go in equal numbers through its two output channels. Classical theory says that the two magnets together have measured both the x - and y -components of spin angular momentum and that the atoms in beam 3 have $S_x = +1/2\hbar$, $S_y = +1/2\hbar$, while those in beam 4 have $S_x = +1/2\hbar$, $S_y = -1/2\hbar$. However, classical theory is wrong, because if beam 3 is put through still another magnet C, with its magnetic field along x , the atoms divide equally into beams 5 and 6 instead of emerging as a single beam 5 (as they would if they had $S_x = +1/2\hbar$). Thus, the correct statement is that the beam entering B has $S_x = +1/2\hbar$ and is composed of an equal mixture of the states $S_y = +1/2\hbar$ and $S_y = -1/2\hbar$ —*i.e.*, the x -component of angular momentum is known but the y -component is not. Correspondingly, beam 3 leaving B has $S_x = +1/2\hbar$ and is an equal mixture of the states $S_y = +1/2\hbar$ and $S_y = -1/2\hbar$; the y -component of angular momentum is known but the x -component is not. The

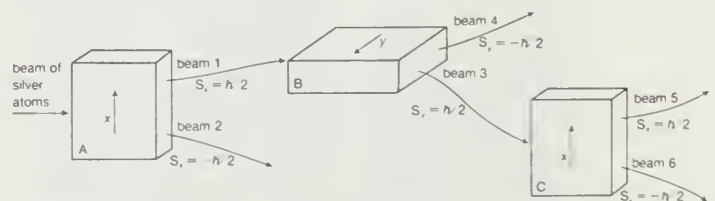


Figure 64: Measurements of the x and y components of angular momentum for silver atoms, S, in the ground state. A, B, and C are magnets with inhomogeneous magnetic fields. The arrows show the average direction of each magnetic field.

Observables with different sets of state functions

information about S_x is lost because of the disturbance caused by magnet B in the measurement of S_y .

HEISENBERG UNCERTAINTY PRINCIPLE

The observables discussed so far have had discrete sets of experimental values. For example, the values of the energy of a bound system are always discrete, and angular momentum components have values that take the form $m\hbar$, where m is either an integer or a half-integer, positive or negative. On the other hand, the position of a particle or the linear momentum of a free particle can take continuous values in both quantum and classical theory. The mathematics of observables with a continuous spectrum of measured values is somewhat more complicated than for the discrete case but presents no problems of principle. An observable with a continuous spectrum of measured values has an infinite number of state functions. The state function Ψ of the system is still regarded as a combination of the state functions of the observable, but the sum in equation (187) must be replaced by an integral.

Measurements can be made of position x of a particle and the x -component of its linear momentum, denoted by p_x . These two observables are incompatible because they have different state functions. The phenomenon of diffraction noted above illustrates the impossibility of measuring position and momentum simultaneously and precisely. If a parallel monochromatic light beam passes through a slit (Figure 65A), its intensity varies with direction, as shown in Figure 65B. The light has zero intensity in certain directions. Wave theory shows that the first zero occurs at an angle θ_0 , given by $\sin \theta_0 = \lambda/b$, where λ is the wavelength of the light and b is the width of the slit. If the width of the slit is reduced, θ_0 increases—*i.e.*, the diffracted light is more spread out. Thus, θ_0 measures the spread of the beam.

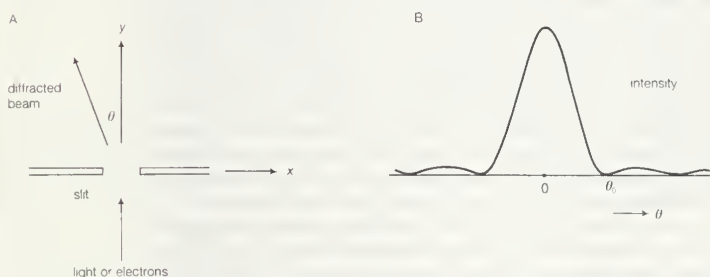


Figure 65: (A) Parallel monochromatic light incident normally on a slit, (B) variation in the intensity of the light with direction after it has passed through the slit. If the experiment is repeated with electrons instead of light, the same diagram would represent the variation in the intensity (*i.e.*, relative number) of the electrons.

The experiment can be repeated with a stream of electrons instead of a beam of light. According to Broglie, electrons have wavelike properties; therefore, the beam of electrons emerging from the slit should widen and spread out like a beam of light waves. This has been observed in experiments. If the electrons have velocity u in the forward direction (*i.e.*, the y -direction in Figure 65A), their (linear) momentum is $p = m_e u$. Consider p_x , the component of momentum in the x -direction. After the electrons have passed through the aperture, the spread in their directions results in an uncertainty in p_x by an amount

$$\Delta p_x \approx p \sin \theta_0 = p \frac{\lambda}{b}, \quad (188)$$

where λ is the wavelength of the electrons and, according to the Broglie formula, equals h/p . Thus, $\Delta p_x \approx h/b$. Exactly where an electron passed through the slit is unknown; it is only certain that an electron went through somewhere. Therefore, immediately after an electron goes through, the uncertainty in its x -position is $\Delta x \approx b/2$. Thus the product of the uncertainties is of the order of \hbar . More exact analysis shows that the product has a lower limit, given by

$$\Delta x \Delta p_x \geq \frac{\hbar}{2}. \quad (189)$$

This is the well-known Heisenberg uncertainty principle for position and momentum. It states that there is a limit to the precision with which the position and the momentum of an object can be measured at the same time. Depending on the experimental conditions, either quantity can be measured as precisely as desired (at least in principle), but the more precisely one of the quantities is measured, the less precisely the other is known.

The uncertainty principle is significant only on the atomic scale because of the small value of h in everyday units. If the position of a macroscopic object with a mass of, say one gram, is measured with a precision of 10^{-6} metre, the uncertainty principle states that its velocity cannot be measured to better than about 10^{-25} metre per second. Such a limitation is hardly worrisome. However, if an electron is located in an atom about 10^{-10} metre across, the principle gives a minimum uncertainty in the velocity of about 10^6 metre per second.

The above reasoning leading to the uncertainty principle is based on the wave-particle duality of the electron. When Heisenberg first propounded the principle in 1927 his reasoning was based, however, on the wave-particle duality of the photon. He considered the process of measuring the position of an electron by observing it in a microscope. Diffraction effects due to the wave nature of light result in a blurring of the image; the resulting uncertainty in the position of the electron is approximately equal to the wavelength of the light. To reduce this uncertainty, it is necessary to use light of shorter wavelength—*e.g.*, gamma rays. However, in producing an image of the electron, the gamma-ray photon bounces off the electron, giving the Compton effect (see above). As a result of the collision, the electron recoils in a statistically random way. The resulting uncertainty in the momentum of the electron is proportional to the momentum of the photon, which is inversely proportional to the wavelength of the photon. So it is again the case that increased precision in knowledge of the position of the electron is gained only at the expense of decreased precision in knowledge of its momentum. A detailed calculation of the process yields the same result as before (equation [189]). Heisenberg's reasoning brings out clearly the fact that the smaller the particle being observed, the more significant is the uncertainty principle. When a large body is observed, photons still bounce off it and change its momentum, but, considered as a fraction of the initial momentum of the body, the change is insignificant.

The Schrödinger and Dirac theories give a precise value for the energy of each stationary state, but in reality the states do not have a precise energy. The only exception is in the ground (lowest energy) state. Instead, the energies of the states are spread over a small range. The spread arises from the fact that, because the electron can make a transition to another state, the initial state has a finite lifetime. The transition is a random process, and so different atoms in the same state have different lifetimes. If the mean lifetime is denoted as τ , the theory shows that the energy of the initial state has a spread of energy ΔE , given by

$$\tau \Delta E \approx \hbar. \quad (190)$$

This energy spread is manifested in a spread in the frequencies of emitted radiation. Therefore, the spectral lines are not infinitely sharp. (Some experimental factors can also broaden a line, but their effects can be reduced; however, the present effect, known as natural broadening, is fundamental and cannot be reduced.) Equation (190) is another type of Heisenberg uncertainty relation; generally, if a measurement with duration τ is made of the energy in a system, the measurement disturbs the system, causing the energy to be uncertain by an amount ΔE , the magnitude of which is given by the above equation.

QUANTUM ELECTRODYNAMICS

The application of quantum theory to the interaction between electrons and radiation requires a quantum treatment of Maxwell's field equations, which are the foundations of electromagnetism, and the relativistic theory of the electron formulated by Dirac (see above *Electron spin and antiparticles*). The resulting quantum field theory is known as quantum electrodynamics, or QED.

Position-momentum uncertainty relation

Energy-time uncertainty principle

QED accounts for the behaviour and interactions of electrons, positrons, and photons. It deals with processes involving the creation of material particles from electromagnetic energy and with the converse processes in which a material particle and its antiparticle annihilate each other and produce energy. Initially the theory was beset with formidable mathematical difficulties, because the calculated values of quantities such as the charge and mass of the electron proved to be infinite. However, an ingenious set of techniques developed (in the late 1940s) by Hans Bethe, Julian S. Schwinger, Tomonaga Shin'ichirō, Richard P. Feynman, and others dealt systematically with the infinities to obtain finite values of the physical quantities. Their method is known as renormalization. The theory has provided some remarkably accurate predictions.

According to the Dirac theory, two particular states in hydrogen with different quantum numbers have the same energy. QED, however, predicts a small difference in their energies; the difference may be determined by measuring the frequency of the electromagnetic radiation that produces transitions between the two states. This effect was first measured by Willis E. Lamb, Jr., and Robert Retherford in 1947. Its physical origin lies in the interaction of the electron with the random fluctuations in the surrounding electromagnetic field. These fluctuations, which exist even in the absence of an applied field, are a quantum phenomenon. The accuracy of experiment and theory in this area may be gauged by two recent values for the separation of the two states, expressed in terms of the frequency of the radiation that produces the transitions:

experiment (1982)	$1,057,858 \pm 2$ kilohertz
theory (1975)	$1,057,864 \pm 14$ kilohertz.

Magnetic dipole moment of the electron

An even more spectacular example of the success of QED is provided by the value for μ_e , the magnetic dipole moment of the free electron. Because the electron is spinning and has electric charge, it behaves like a tiny magnet, the strength of which is expressed by the value of μ_e . According to the Dirac theory, μ_e is exactly equal to $\mu_B = e\hbar/2m_e$, a quantity known as the Bohr magneton; however, QED predicts that $\mu_e = (1 + a)\mu_B$, where a is a small number, approximately $1/860$. Again, the physical origin of the QED correction is the interaction of the electron with random oscillations in the surrounding electromagnetic field. The best experimental determination of μ_e involves measuring not the quantity itself but the small correction term $\mu_e - \mu_B$. This greatly enhances the sensitivity of the experiment. The most recent results for the value of a are

experiment (1984)	$(115,965,219 \pm 1) \times 10^{-11}$
theory (1986)	$(115,965,227 \pm 10) \times 10^{-11}$.

Since a itself represents a small correction term, the magnetic dipole moment of the electron is measured with an accuracy of about one part in 10^{11} . One of the most precisely determined quantities in physics, the magnetic dipole moment of the electron can be calculated correctly from quantum theory to within about one part in 10^{10} .

The interpretation of quantum mechanics

Although quantum mechanics has been applied to problems in physics with great success, some of its ideas seem strange. A few of their implications are considered here.

THE ELECTRON: WAVE OR PARTICLE?

Young's aforementioned experiment in which a parallel beam of monochromatic light is passed through a pair of narrow parallel slits (Figure 66A) has an electron counterpart. In Young's original experiment, the intensity of the light varies with direction after passing through the slits (Figure 66B). The intensity oscillates because of interference between the light waves emerging from the two slits, the rate of oscillation depending on the wavelength of the light and the separation of the slits. The oscillation creates a fringe pattern of alternating light and dark bands that is modulated by the diffraction pattern from each slit. If one of the slits is covered, the interference fringes disappear, and only the diffraction pattern (shown as a broken line in Figure 66B) is observed.

Two-slit experiment

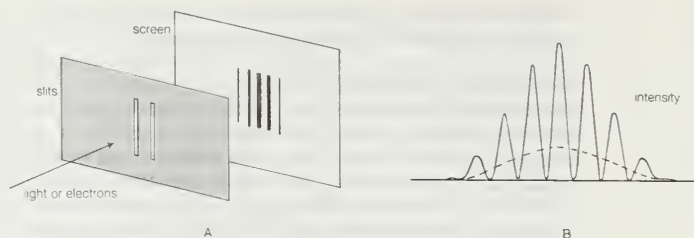


Figure 66: (A) Monochromatic light incident on a pair of slits gives interference fringes (alternate light and dark bands) on a screen. (B) variation in the intensity of the light at the screen when both slits are open. With a single slit, there is no interference pattern; the intensity variation is shown by the broken line. As with Figure 65B, the same diagram would give the variation in the intensity of electrons in the corresponding electron experiment.

Young's experiment can be repeated with electrons all with the same momentum. The screen in the optical experiment is replaced by a closely spaced grid of electron detectors. There are many devices for detecting electrons; the most common are scintillators. When an electron passes through a scintillating material, such as sodium iodide, the material produces a light flash which gives a voltage pulse that can be amplified and recorded. The pattern of electrons recorded by each detector is the same as that predicted for waves with wavelengths given by the Broglie formula. Thus, the experiment provides conclusive evidence for the wave behaviour of electrons.

If the experiment is repeated with a very weak source of electrons so that only one electron passes through the slits, a single detector registers the arrival of an electron. This is a well-localized event characteristic of a particle. Each time the experiment is repeated, one electron passes through the slits and is detected. A graph plotted with detector position along one axis and the number of electrons along the other looks exactly like the oscillating interference pattern in Figure 66B. Thus, the intensity function in the figure is proportional to the probability of the electron moving in a particular direction after it has passed through the slits. Apart from its units, the function is identical to Ψ^2 , where Ψ is the solution of the time-independent Schrödinger equation for this particular experiment.

If one of the slits is covered, the fringe pattern disappears and is replaced by the diffraction pattern for a single slit. Thus, both slits are needed to produce the fringe pattern. However, if the electron is a particle, it seems reasonable to suppose that it passed through only one of the slits. The apparatus can be modified to ascertain which slit by placing a thin wire loop around each slit. When an electron passes through a loop, it generates a small electric signal, showing which slit it passed through. However, the interference fringe pattern then disappears, and the single-slit diffraction pattern returns. Since both slits are needed for the interference pattern to appear and since it is impossible to know which slit the electron passed through without destroying that pattern, one is forced to the conclusion that the electron goes through both slits at the same time.

In summary, the experiment shows both the wave and particle properties of the electron. The wave property predicts the probability of direction of travel before the electron is detected; on the other hand, the fact that the electron is detected in a particular place shows that it has particle properties. Therefore, the answer to the question whether the electron is a wave or a particle is that it is neither. It is an object exhibiting either wave or particle properties, depending on the type of measurement that is made on it. In other words, one cannot talk about the intrinsic properties of an electron; instead, one must consider the properties of the electron and measuring apparatus together.

HIDDEN VARIABLES

A fundamental concept in quantum mechanics is that of randomness, or indeterminacy. In general, the theory predicts only the probability of a certain result. Consider the case of radioactivity. Imagine a box of atoms with identi-

Experimental proof of wave-particle duality

The indeterminacy of quantum mechanics

cal nuclei that can undergo decay with the emission of an alpha particle. In a given time interval, a certain fraction will decay. The theory may tell precisely what that fraction will be, but it cannot predict which particular nuclei will decay. The theory asserts that, at the beginning of the time interval, all the nuclei are in an identical state and that the decay is a completely random process. Even in classical physics, many processes appear random. For example, one says that, when a roulette wheel is spun, the ball will drop at random into one of the numbered compartments in the wheel. Based on this belief, the casino owner and the players give and accept identical odds against each number for each throw. However, the fact is that the winning number could be predicted if one noted the exact location of the wheel when the croupier released the ball, the initial speed of the wheel, and various other physical parameters. It is only ignorance of the initial conditions and the difficulty of doing the calculations that makes the outcome appear to be random. In quantum mechanics, on the other hand, the randomness is asserted to be absolutely fundamental. The theory says that, though one nucleus decayed and the other did not, they were previously in the identical state.

Many eminent physicists, including Einstein, have not accepted this indeterminacy. They have rejected the notion that the nuclei were initially in the identical state. Instead, they postulated that there must be some other property—presently unknown, but existing nonetheless—that is different for the two nuclei. This type of unknown property is termed a hidden variable; if it existed, it would restore determinacy to physics. If the initial values of the hidden variables were known, it would be possible to predict which nuclei would decay. Such a theory would, of course, also have to account for the wealth of experimental data which conventional quantum mechanics explains from a few simple assumptions. Attempts have been made by Broglie, David Bohm, and others to construct theories based on hidden variables, but the theories are very complicated and contrived. For example, the electron would definitely have to go through only one slit in the two-slit experiment. To explain that interference occurs only when the other slit is open, it is necessary to postulate a special force on the electron which exists only when that slit is open. Such artificial additions make hidden variable theories unattractive, and there is little support for them among physicists.

The orthodox view of quantum mechanics—and the one adopted in the present article—is known as the Copenhagen interpretation because its main protagonist, Niels Bohr, worked in that city. The Copenhagen view of understanding the physical world stresses the importance of basing theory on what can be observed and measured experimentally. It therefore rejects the idea of hidden variables as quantities that cannot be measured. The Copenhagen view is that the indeterminacy observed in nature is fundamental and does not reflect an inadequacy in present scientific knowledge. One should therefore accept the indeterminacy without trying to “explain” it and see what consequences come from it.

Attempts have been made to link the existence of free will with the indeterminacy of quantum mechanics, but it is difficult to see how this feature of the theory makes free will more plausible. On the contrary, free will presumably implies rational thought and decision, whereas the essence of the indeterminism in quantum mechanics is that it is due to intrinsic randomness.

PARADOX OF EINSTEIN, PODOLSKY, AND ROSEN

In 1935 Einstein and two other physicists in the United States, Boris Podolsky and Nathan Rosen, analyzed a thought experiment to measure position and momentum in a pair of interacting systems. Employing conventional quantum mechanics, they obtained some startling results, which led them to conclude that the theory does not give a complete description of physical reality. Their results, which are so peculiar as to seem paradoxical, are based on impeccable reasoning, but their conclusion that the theory is incomplete does not necessarily follow. Bohm simplified their experiment while retaining the central point of their reasoning; this discussion follows his account.

The proton, like the electron, has spin $1/2$; thus, no matter what direction is chosen for measuring the component of its spin angular momentum, the values are always $+1/2\hbar$ or $-1/2\hbar$. (The present discussion relates only to spin angular momentum, and the word spin is omitted from now on.) It is possible to obtain a system consisting of a pair of protons in close proximity and with total angular momentum equal to zero. Thus, if the value of one of the components of angular momentum for one of the protons is $+1/2\hbar$ along any selected direction, the value for the component in the same direction for the other particle must be $-1/2\hbar$. Suppose the two protons move in opposite directions until they are far apart. The total angular momentum of the system remains zero, and if the component of angular momentum along the same direction for each of the two particles is measured, the result is a pair of equal and opposite values. Therefore, after the quantity is measured for one of the protons, it can be predicted for the other proton; the second measurement is unnecessary. As previously noted, measuring a quantity changes the state of the system. Thus, if measuring S_x (the x -component of angular momentum) for proton 1 produces the value $+1/2\hbar$, the state of proton 1 after measurement corresponds to $S_x = +1/2\hbar$, and the state of proton 2 corresponds to $S_x = -1/2\hbar$. Any direction, however, can be chosen for measuring the component of angular momentum. Whichever direction is selected, the state of proton 1 after measurement corresponds to a definite component of angular momentum about that direction. Furthermore, since proton 2 must have the opposite value for the same component, it follows that the measurement on proton 1 results in a definite state for proton 2 relative to the chosen direction, notwithstanding the fact that the two particles may be millions of kilometres apart and are not interacting with each other at the time. Einstein and his two collaborators thought that this conclusion was so obviously false that the quantum mechanical theory on which it was based must be incomplete. They concluded that the correct theory would contain some hidden variable feature that would restore the determinism of classical physics.

A comparison of how quantum theory and classical theory describe angular momentum for particle pairs illustrates the essential difference between the two outlooks. In both theories, if a system of two particles has a total angular momentum of zero, then the angular momenta of the two particles are equal and opposite. If the components of angular momentum are measured along the same direction, the two values are numerically equal, one positive and the other negative. Thus, if one component is measured, the other can be predicted. The crucial difference between the two theories is that, in classical physics, the system under investigation is assumed to have possessed the quantity being measured beforehand. The measurement does not disturb the system; it merely reveals the preexisting state. It may be noted that, if a particle were actually to possess components of angular momentum prior to measurement, such quantities would constitute hidden variables.

Does nature behave as quantum mechanics predicts? The answer comes from measuring the components of angular momenta for the two protons along different directions with an angle θ between them. A measurement on one proton can give only the result $+1/2\hbar$ or $-1/2\hbar$. The experiment consists of measuring correlations between the plus and minus values for pairs of protons with a fixed value of θ , and then repeating the measurements for different values of θ , as in Figure 67. The interpretation of the results rests on an important theorem by the British physicist John Stewart Bell. Bell began by assuming the existence of some form of hidden variable with a value that would determine whether the measured angular momentum gives a plus or minus result. He further assumed locality—namely, that measurement on one proton (*i.e.*, the choice of the measurement direction) cannot affect the result of the measurement on the other proton. Both these assumptions agree with classical, commonsense ideas. He then showed quite generally that these two assumptions lead to a certain relationship, now known as Bell's inequality, for the correlation values mentioned above. Experiments

Essential difference between quantum mechanics and classical physics

Rejection of the notion of hidden variables

Bell's inequality

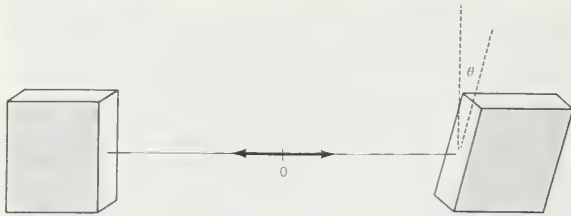


Figure 67: Experiment to determine the correlation in measured angular momentum values for a pair of protons with zero total angular momentum. The two protons are initially at the point 0 and move in opposite directions toward the two magnets.

have been conducted at several laboratories with photons instead of protons (the analysis is similar), and the results show fairly conclusively that Bell's inequality is violated. That is to say, the observed results agree with those of quantum mechanics and cannot be accounted for by a hidden variable (or deterministic) theory based on the concept of locality. One is forced to conclude that the two protons are a correlated pair and that a measurement on one affects the state of both, no matter how far apart they are. This may strike one as highly peculiar, but such is the way nature appears to be.

It may be noted that the effect on the state of proton 2 following a measurement on proton 1 is believed to be instantaneous; the effect happens before a light signal initiated by the measuring event at proton 1 reaches proton 2. Alain Aspect and his coworkers in Paris demonstrated this result in 1982 with an ingenious experiment in which the correlation between the two angular momenta was measured, within a very short time interval, by a high-frequency switching device. The interval was less than the time taken for a light signal to travel from one particle to the other at the two measurement positions. Einstein's special theory of relativity states that no message can travel with a speed greater than that of light. Thus, there is no way that the information concerning the direction of the measurement on the first proton could reach the second proton before the measurement was made on it.

MEASUREMENT IN QUANTUM MECHANICS

The way quantum mechanics treats the process of measurement has caused considerable debate. Schrödinger's time-dependent wave equation (equation [185]) is an exact recipe for determining the way the wave function varies with time for a given physical system in a given physical environment. According to the Schrödinger equation, the wave function varies in a strictly determinate way. On the other hand, in the axiomatic approach to quantum mechanics described above, a measurement changes the wave function abruptly and discontinuously. Before the measurement is made, the wave function Ψ is a mixture of the ϕ s as indicated in equation (187). The measurement changes Ψ from a mixture of ϕ s to a single ϕ . This change, brought about by the process of measurement, is termed the collapse or reduction of the wave function. The collapse is a discontinuous change in Ψ ; it is also unpredictable, because, starting with the same Ψ represented by the right-hand side of equation (187), the end result can be any one of the individual ϕ s.

The Schrödinger equation, which gives a smooth and predictable variation of Ψ , applies between the measurements. The measurement process itself, however, cannot be described by the Schrödinger equation; it is somehow a thing apart. This appears unsatisfactory, inasmuch as a measurement is a physical process and ought to be the subject of the Schrödinger equation just like any other physical process.

The difficulty is related to the fact that quantum mechanics applies to microscopic systems containing one (or a few) electrons, protons, or photons. Measurements, however, are made with large-scale objects (*e.g.*, detectors, amplifiers, and meters) in the macroscopic world, which obeys the laws of classical physics. Thus, another way of formulating the question of what happens in a measurement is to ask how the microscopic quantum world relates and interacts with the macroscopic classical world. More

narrowly, it can be asked how and at what point in the measurement process does the wave function collapse? So far, there are no satisfactory answers to these questions, although there are several schools of thought.

One approach stresses the role of a conscious observer in the measurement process and suggests that the wave function collapses when the observer reads the measuring instrument. Bringing the conscious mind into the measurement problem seems to raise more questions than it answers, however.

As discussed above, the Copenhagen interpretation of the measurement process is essentially pragmatic. It distinguishes between microscopic quantum systems and macroscopic measuring instruments. The initial object or event—*e.g.*, the passage of an electron, photon, or atom—triggers the classical measuring device into giving a reading; somewhere along the chain of events, the result of the measurement becomes fixed (*i.e.*, the wave function collapses). This does not answer the basic question but says, in effect, not to worry about it. This is probably the view of most practicing physicists.

A third school of thought notes that an essential feature of the measuring process is irreversibility. This contrasts with the behaviour of the wave function when it varies according to the Schrödinger equation; in principle, any such variation in the wave function can be reversed by an appropriate experimental arrangement. However, once a classical measuring instrument has given a reading, the process is not reversible. It is possible that the key to the nature of the measurement process lies somewhere here. The Schrödinger equation is known to apply only to relatively simple systems. It is an enormous extrapolation to assume that the same equation applies to the large and complex system of a classical measuring device. It may be that the appropriate equation for such a system has features that produce irreversible effects (*e.g.*, wave-function collapse) which differ in kind from those for a simple system.

One may also mention the so-called many-worlds interpretation, proposed by Hugh Everett III in 1957, which suggests that, when a measurement is made for a system in which the wave function is a mixture of states, the universe branches into a number of noninteracting universes. Each of the possible outcomes of the measurement occurs, but in a different universe. Thus, if $S_x = \hbar/2$ is the result of a Stern-Gerlach measurement on a silver atom (see above), there is another universe identical to ours in every way (including clones of people), except that the result of the measurement is $S_x = -\hbar/2$. Although this fanciful model solves some measurement problems, it has few adherents among physicists.

Because the various ways of looking at the measurement process lead to the same experimental consequences, trying to distinguish between them on scientific grounds may be fruitless. One or another may be preferred on the grounds of plausibility, elegance, or economy of hypotheses, but these are matters of individual taste. Whether one day a satisfactory quantum theory of measurement will emerge, distinguished from the others by its verifiable predictions, remains an open question.

Applications of quantum mechanics

As has been noted, quantum mechanics has been enormously successful in explaining microscopic phenomena in all branches of physics. The three phenomena described in this section are examples that demonstrate the quintessence of the theory.

DECAY OF THE K^0 MESON

The K^0 meson, discovered in 1953, is produced in high-energy collisions between nuclei and other particles. It has zero electric charge, and its mass is about one-half the mass of the proton. It is unstable and, once formed, rapidly decays into either 2 or 3 pi-mesons. The average lifetime of the K^0 is about 10^{-10} second.

In spite of the fact that the K^0 meson is uncharged, quantum theory predicts the existence of an antiparticle with the same mass, decay products, and average lifetime; the

Copen-
hagen
interpreta-
tion

The many-
worlds
interpreta-
tion

Collapse of
the wave
function Ψ

antiparticle is denoted by \bar{K}^0 . During the early 1950s, several physicists questioned the justification for postulating the existence of two particles with such similar properties. In 1955, however, Murray Gell-Mann and Abraham Pais made an interesting prediction about the decay of the K^0 meson. Their reasoning provides an excellent illustration of the quantum mechanical axiom that the wave function Ψ can be a superposition of states; in this case, there are two states, the K^0 and \bar{K}^0 mesons themselves.

A K^0 meson may be represented formally by writing the wave function as $\Psi = K^0$; similarly $\Psi = \bar{K}^0$ represents a \bar{K}^0 . From the two states, K^0 and \bar{K}^0 , the following two new states are constructed:

$$K_1 = \frac{(K^0 + \bar{K}^0)}{\sqrt{2}}, \tag{191a}$$

$$K_2 = \frac{(K^0 - \bar{K}^0)}{\sqrt{2}}. \tag{191b}$$

From these two equations it follows that

$$K^0 = \frac{(K_1 + K_2)}{\sqrt{2}}, \tag{191c}$$

$$\bar{K}^0 = \frac{(K_1 - K_2)}{\sqrt{2}}. \tag{191d}$$

The reason for defining the two states K_1 and K_2 is that, according to quantum theory, when the K^0 decays, it does not do so as an isolated particle; instead, it combines with its antiparticle to form the states K_1 and K_2 . The state K_1 decays into two pi-mesons with a very short lifetime (about 10^{-10} second), while K_2 decays into three pi-mesons with a longer lifetime (about 10^{-7} second).



Figure 68: Decay of the K^0 meson.

The physical consequences of these results may be demonstrated in the following experiment. K^0 particles are produced in a nuclear reaction at the point A (Figure 68). They move to the right in the figure and start to decay. At point A, the wave function is $\Psi = K^0$, which, from equation (191c), can be expressed as the sum of K_1 and K_2 . As the particles move to the right, the K_1 state begins to decay rapidly. If the particles reach point B in about 10^{-8} second, nearly all the K_1 component has decayed, although hardly any of the K_2 component has done so. Thus, at point B, the beam has changed from one of pure K^0 to one of almost pure K_2 , which equation (191b) shows is an equal mixture of K^0 and \bar{K}^0 . In other words, \bar{K}^0 particles appear in the beam simply because K_1 and K_2 decay at different rates. At point B, the beam enters a block of absorbing material. Both the K^0 and \bar{K}^0 are absorbed by the nuclei in the block, but the \bar{K}^0 are absorbed more strongly. As a result, even though the beam is an equal mixture of K^0 and \bar{K}^0 when it enters the absorber, it is almost pure K^0 when it exits at point C. The beam thus begins and ends as K^0 .

Gell-Mann and Pais predicted all this, and experiments subsequently verified it. The experimental observations are that the decay products are primarily two pi-mesons with a short decay time near A, three pi-mesons with longer decay time near B, and two pi-mesons again near C. (This account exaggerates the changes in the K_1 and K_2 components between A and B and in the K^0 and \bar{K}^0 components between B and C; the argument, however, is unchanged.) The phenomenon of generating the \bar{K}^0 and regenerating the K_1 decay is purely quantum. It rests on the quantum axiom of the superposition of states and has no classical counterpart.

CESIUM CLOCK

The cesium clock is the most accurate type of clock yet developed. This device makes use of transitions between the spin states of the cesium nucleus and produces a frequency which is so regular that it has been adopted for establishing the time standard.

Quantum oscillation frequency

Like electrons, many atomic nuclei have spin. The spin of these nuclei produces a set of small effects in the spectra, known as hyperfine structure. (The effects are small because, though the angular momentum of a spinning nucleus is of the same magnitude as that of an electron, its magnetic moment, which governs the energies of the atomic levels, is relatively small.) The nucleus of the cesium atom has spin quantum number $7/2$. The total angular momentum of the lowest energy states of the cesium atom is obtained by combining the spin angular momentum of the nucleus with that of the single valence electron in the atom. (Only the valence electron contributes to the angular momentum because the angular momenta of all the other electrons total zero. Another simplifying feature is that the ground states have zero orbital momenta, so only spin angular momenta need to be considered.) When nuclear spin is taken into account, the total angular momentum of the atom is characterized by a quantum number, conventionally denoted by F , which for cesium is 4 or 3. These values come from the spin value $7/2$ for the nucleus and $1/2$ for the electron. If the nucleus and the electron are visualized as tiny spinning tops, the value $F = 4$ ($7/2 + 1/2$) corresponds to the tops spinning in the same sense, and $F = 3$ ($7/2 - 1/2$) corresponds to spins in opposite senses. The energy difference ΔE of the states with the two F values is a precise quantity. If electromagnetic radiation of frequency ν_0 , where

$$h\nu_0 = \Delta E, \tag{192}$$

is applied to a system of cesium atoms, transitions will occur between the two states. An apparatus that can detect the occurrence of transitions thus provides an extremely precise frequency standard. This is the principle of the cesium clock.

Principle of the cesium clock

The apparatus is shown schematically in Figure 69. A beam of cesium atoms emerges from an oven at a temperature of about 100° C. The atoms pass through an inhomogeneous magnet A, which deflects the atoms in state $F = 4$ downward and those in state $F = 3$ by an equal amount upward. The atoms pass through slit S and continue into a second inhomogeneous magnet B. Magnet B is arranged so that it deflects atoms with an unchanged state in the same direction that magnet A deflected them. The atoms follow the paths indicated by the broken lines in the figure and are lost to the beam. However, if an alternating electromagnetic field of frequency ν_0 is applied to the beam as it traverses the centre region C, transitions between states will occur. Some atoms in state $F = 4$ will change to $F = 3$, and vice versa. For such atoms, the deflections in magnet B are reversed. The atoms follow the whole lines in the diagram and strike a tungsten wire, which gives electric signals in proportion to the number of cesium atoms striking the wire. As the frequency ν of the alternating field is varied, the signal has a sharp maximum for $\nu = \nu_0$. The length of the apparatus from

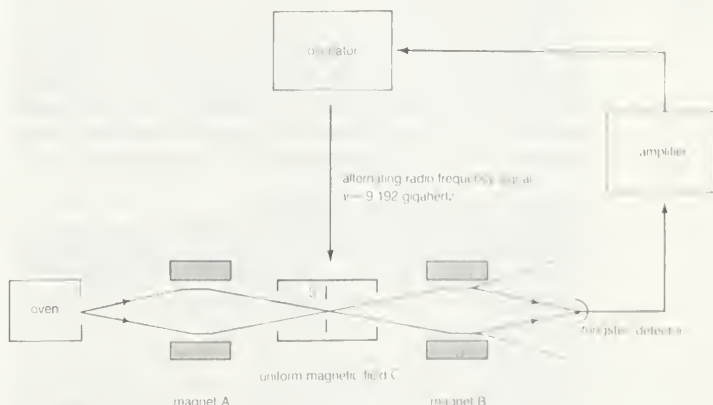


Figure 69: Cesium clock.

the oven to the tungsten detector is about one metre.

Each atomic state is characterized not only by the quantum number F but also by a second quantum number m_F . For $F=4$, m_F can take integral values from 4 to -4 . In the absence of a magnetic field, these states have the same energy. A magnetic field, however, causes a small change in energy proportional to the magnitude of the field and to the m_F value. Similarly, a magnetic field changes the energy for the $F=3$ states according to the m_F value which, in this case, may vary from 3 to -3 . The energy changes are indicated in Figure 70. In the cesium clock, a weak constant magnetic field is superposed on the alternating electromagnetic field in region C. The theory shows that the alternating field can bring about a transition only between pairs of states with m_F values that are the same or that differ by unity. However, as can be seen from the figure, the only transitions occurring at the frequency ν_0 are those between the two states with $m_F=0$. The apparatus is so sensitive that it can discriminate easily between such transitions and all the others.

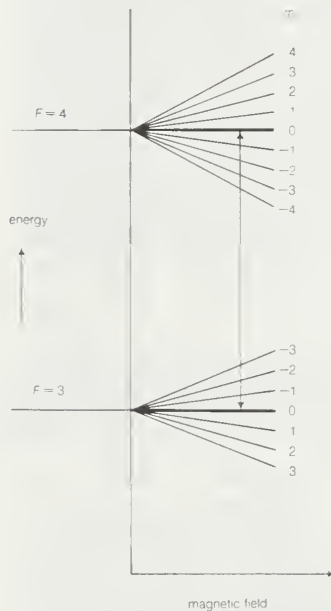


Figure 70: Variation of energy with magnetic-field strength for the $F=4$ and $F=3$ states in cesium-133.

If the frequency of the oscillator drifts slightly so that it does not quite equal ν_0 , the detector output drops. The change in signal strength produces a signal to the oscillator to bring the frequency back to the correct value. This servo system keeps the oscillator frequency automatically locked to ν_0 .

The cesium clock is exceedingly stable. The frequency of the oscillator remains constant to about one part in 10^{13} . For this reason, the device has been used to redefine the second. This base unit of time in the SI system is defined as equal to 9,192,631,770 cycles of the radiation corresponding to the transition between the levels $F=4$, $m_F=0$ and $F=3$, $m_F=0$ of the ground state of the cesium-133 atom. Prior to 1964, the second was defined in terms of the motion of the Earth. The latter, however, is not nearly as stable as the cesium clock. Specifically, the fractional variation of the Earth's rotation is a few hundred times larger than that of the frequency of the cesium clock.

A QUANTUM VOLTAGE STANDARD

Quantum theory has been used to establish a voltage standard, and this standard has proven to be extraordinarily accurate and consistent from laboratory to laboratory.

If two layers of superconducting material are separated by a thin insulating barrier, a supercurrent (*i.e.*, a current of paired electrons) can pass from one superconductor to the other. This is another example of the tunneling process described earlier. Several effects based on this phenomenon were predicted in 1962 by the British physicist Brian D.

Josephson. Demonstrated experimentally soon afterwards, they are now referred to as the Josephson effects.

If a DC (direct-current) voltage V is applied across the two superconductors, the energy of an electron pair changes by an amount of $2eV$ as it crosses the junction. As a result, the supercurrent oscillates with frequency ν given by the Planck relationship—equation (176). Thus,

$$2eV = h\nu. \quad (193)$$

This oscillatory behaviour of the supercurrent is known as the AC (alternating-current) Josephson effect. Measurement of V and ν permits a direct verification of the Planck relationship. Although the oscillating supercurrent has been detected directly, it is extremely weak. A more sensitive method of investigating equation (193) is to study effects resulting from the interaction of microwave radiation with the supercurrent.

Several carefully conducted experiments have verified equation (193) to such a high degree of precision that it has been used to determine the value of $2e/h$. This value can in fact be determined more precisely by the AC Josephson effect than by any other method. The result is so reliable that laboratories now employ the AC Josephson effect to set a voltage standard. The numerical relationship between V and ν is

$$\frac{2e}{h} = \frac{\nu}{V} = 483,597.7 \times 10^9 \text{ hertz per volt.}$$

In this way, measuring a frequency, which can be done with great precision, gives the value of the voltage. Before the Josephson method was used, the voltage standard in metrological laboratories devoted to the maintenance of physical units was based on high-stability Weston cadmium cells. These cells, however, tend to drift and so caused inconsistencies between standards in different laboratories. The Josephson method has provided a standard giving agreement to within a few parts in 10^8 for measurements made at different times and in different laboratories.

The experiments described in the preceding two sections are only two examples of high-precision measurements in physics. The values of the fundamental constants, such as c , h , e , and m_e , are determined from a wide variety of experiments based on quantum phenomena. The results are so consistent that the values of the constants are thought to be known in most cases to better than one part in 10^6 . Physicists may not know what they are doing when they make a measurement, but they do it extremely well. (G.L.S.)

BIBLIOGRAPHY

Classical mechanics. The history of classical mechanics is chronicled in I. BERNARD COHEN, *The Newtonian Revolution* (1980); and E.J. DIJKSTERHUIS, *The Mechanization of the World Picture* (1961; originally published in Dutch, 1950). All introductory physics textbooks contain a portion on classical mechanics; recent examples include HANS C. OGANIAN, *Physics*, 2nd ed. (1989); and ROBERT RESNICK, DAVID HALLIDAY, and KENNETH S. KRANE, *Physics*, 4th ed. (1992). The principal reference for classical mechanics is S. ERAUTSCHI *et al.*, *The Mechanical Universe: Mechanics and Heat*, advanced ed. (1986). Other texts include, at an introductory level, A.P. FRENCH, *Newtonian Mechanics* (1971); and at a more advanced, graduate-school level, HERBERT GOLDSTEIN, *Classical Mechanics*, 2nd ed. (1980), a standard text that contains a lengthy, detailed bibliography of more specialized books dealing with specific aspects of the subject. (D.L.G.)

Celestial mechanics. Modern introductory treatments and discussions of some advanced techniques and classic developments include J.M.A. DANBY, *Fundamentals of Celestial Mechanics*, 2nd ed., rev. and enlarged (1988); DIRK BROUWER and GERALD M. CLEMENCE, *Methods of Celestial Mechanics* (1961); and HENRY CROZIER KEATING PLUMMER, *An Introductory Treatise on Dynamical Astronomy* (1918, reprinted 1960). Orbital resonances are discussed in two review articles by S.J. PEALE: "Orbital Resonances in Solar-System," *Annual Review of Astronomy and Astrophysics*, 14:215–246 (1976), and "Orbital Resonances, Unusual Configurations, and Exotic Rotation States Among Planetary Satellites," in JOSEPH A. BURNS and MILDRED SHAPLEY MATTHEWS (eds.), *Satellites* (1986), pp. 159–223. Current practice in solving the n -body problem on computers is given in the introduction to a paper by LARS HERNQUIST, "Performance Characteristics of Tree Codes," *The*

The AC Josephson effect

Definition of the second

Astrophysical Journal: Supplement Series, 64(4):715–734 (August 1987). An introduction to modern dynamics involving chaos and an introduction to algebraic maps is given by MICHEL HENON, "Numerical Exploration of Hamiltonian Systems," in GÉRARD IOOSS, ROBERT H.G. HELLEMAN, and RAYMOND STORA (eds.), *Chaotic Behaviour in Deterministic Systems* (1983), pp. 54–170. Readable accounts of examples of chaotic dynamics in celestial mechanics are found in two articles by JACK WISDOM: "Chaotic Dynamics in the Solar-System," *Icarus*, 72(2):241–275 (1987), and "Chaotic Behaviour in the Solar System," in M.V. BERRY, I.C. PERCIVAL, and N.O. WEISS (eds.), *Dynamical Chaos* (1987), pp. 109–129. A simple discussion of tides and tidal evolution is given by S.J. PEALE, "Consequences of Tidal Evolution," in MARGARET G. KIVELSON (ed.), *The Solar System: Observations and Interpretations* (1986), pp. 275–288. Advanced discussions of tidal evolution analysis as applied to the Earth are given by KURT LAMBECK, *The Earth's Variable Rotation: Geophysical Causes and Consequences* (1980). (S.J.P.)

Relativistic mechanics. An outstanding work containing an account of the special theory of relativity is ABRAHAM PAIS, "Subtle Is the Lord—": *The Science and Life of Albert Einstein* (1982). Some good introductions at the undergraduate level are W. RINDLER, *Essential Relativity: Special, General, and Cosmological*, 2nd ed. (1977); JAMES H. SMITH, *Introduction to Special Relativity* (1965); EDWIN F. TAYLOR and JOHN ARCHIBALD WHEELER, *Spacetime Physics* (1966). More substantial treatises are J. AHARONI, *The Special Theory of Relativity*, 2nd ed. (1965, reprinted 1985); and J.L. SYNGE, *Relativity: The Special Theory*, 2nd ed. (1965). (G.W.G.)

Mechanics of solids. There are a number of works on the history of the subject. A.E.H. LOVE, *A Treatise on the Mathematical Theory of Elasticity*, 4th ed. (1927, reprinted 1944), has a well-researched chapter on the origin of elasticity up to the early 1900s. STEPHEN P. TIMOSHENKO, *History of Strength of Materials: With a Brief Account of the History of Theory of Elasticity and Theory of Structures* (1953, reprinted 1983), provides good coverage of most subfields of solid mechanics up to the period around 1940, including in some cases detailed but quite readable accounts of specific developments and capsule biographies of major figures. C. TRUESDELL, *Essays in the History of Mechanics* (1968), summarizes his studies of original source materials on Jakob Bernoulli (1654–1705), Leonhard Euler, Leonardo da Vinci, and others and connects those contributions to some of the developments in what he calls "rational mechanics" as of the middle 1900s. Two articles in *Handbuch der Physik* provide historical background: C. TRUESDELL and R.A. TOUPIN, "The Classical Field Theories," vol. 3, pt. 1 (1960); and C. TRUESDELL and W. NOLL, "The Nonlinear Field Theories of Mechanics," vol. 3, pt. 3 (1965).

There are many good books for beginners on the subject, intended for the education of engineers; one that stands out for its coverage of inelastic solid mechanics as well as the more conventional topics on elementary elasticity and structures is STEPHEN H. CRANDALL, NORMAN C. DAHL, and THOMAS J. LARDNER (eds.), *An Introduction to the Mechanics of Solids*, 2nd ed., with SI units (1978). Those with an interest in the physics of materials might begin with A.H. COTTRELL, *The Mechanical Properties of Matter* (1964, reprinted 1981). Some books for beginners aim for a more general introduction to continuum mechanics, including solids and fluids; one such text is Y.C. FUNG, *A First Course in Continuum Mechanics*, 2nd ed. (1977). A readable introduction to continuum mechanics at a more advanced level, such as might be used by scientists and engineers from other fields or by first-year graduate students, is LAWRENCE E. MALVERN, *Introduction to the Mechanics of a Continuous Medium* (1969). The article by TRUESDELL and TOUPIN, mentioned above, provides a comprehensive, perhaps overwhelming, treatment of continuum mechanics fundamentals.

For more specialized treatment of linear elasticity, the classics are the work by LOVE, mentioned above; STEPHEN P. TIMOSHENKO and J.N. GOODIER, *Theory of Elasticity*, 3rd ed. (1970); and N.I. MUSKHELISHVILI, *Some Basic Problems of the Mathematical Theory of Elasticity*, 2nd ed. (1963, reprinted 1977; originally published in Russian, 4th corrected and augmented ed., 1954). The article by TRUESDELL and NOLL noted above is a good source on finite elasticity and also on viscoelastic fluids; a standard reference on the latter is R. BYRON BIRD *et al.*, *Dynamics of Polymeric Liquids*, vol. 1, *Fluid Mechanics*, 2nd ed. (1987). Other books generally regarded as classics in their subfields are R. HILL, *The Mathematical Theory of Plasticity* (1950, reissued 1983); J.C. JAEGER and N.G. COOK, *Fundamen-*

tals of Rock Mechanics, 3rd ed. (1979). JOHN PRICE HIRTH and JENS LOTHE, *Theory of Dislocations*, 2nd ed. (1982); and KEIICHI AKI and PAUL G. RICHARDS, *Quantitative Seismology*, 2 vol. (1980). Other aspects of stress waves in solids are covered by J.D. ACHENBACH, *Wave Propagation in Elastic Solids* (1973). In addition, the scope of finite element analysis in solid mechanics and many other areas can be gleaned from o.c. ZIENKIEWICZ and R.L. TAYLOR, *The Finite Element Method*, 4th ed., 2 vol. (1989–91); and that of fracture mechanics from MELVIN F. KANNINEN and CARL H. POPELAR, *Advanced Fracture Mechanics* (1985). Structural mechanics and issues relating to stability and elastic-plastic stress-strain relations in a way that updates the book by Hill are presented by ZDEŇEK P. BĀZANT and LUIGI CEDOLIN, *Stability of Structures: Elastic, Inelastic, Fracture, and Damage Theories* (1991). (J.R.R.)

Fluid mechanics. A classic text which enshrines all the results of 19th-century fluid dynamics is HORACE LAMB, *Hydrodynamics*, 6th ed. (1932, reissued 1945). This remains useful, but many later books, besides being more up-to-date, provide a better-balanced perspective of the subject and have better illustrations. N. CURLE and H.J. DAVIES, *Modern Fluid Dynamics*, vol. 1, *Incompressible Flow* (1968); and G.K. BATCHELOR, *Introduction to Fluid Dynamics* (1967, reissued 1973), can both be recommended to serious students who are not put off by mathematics. D.J. TRITTON, *Physical Fluid Dynamics*, 2nd ed. (1988), adopts a somewhat different approach and contains interesting material on turbulence and convective instabilities. Readers who are interested in the practical aspects of the subject and who want information concerning hydrostatics as well as fluid dynamics should consult one of the many good texts intended for engineers; among these, B.S. MASSEY, *Mechanics of Fluids*, 5th ed. (1983), is excellent. The development of the subject as a practical science is traced in HUNTER ROUSE and SIMON INCE, *History of Hydraulics* (1957, reprinted 1980). (T.E.F.)

Quantum mechanics. Scholarly but lively accounts of the early history of the subject are JAGDISH MEHRA and HELMUT RECHENBERG, *The Historical Development of Quantum Theory* (1982–), of which 5 vol. in 7 had appeared by 1992, covering the period from 1900 to 1926; and MAX JAMMER, *The Conceptual Development of Quantum Mechanics* (1966). The authors had the advantage, not usually accorded to historians, of being able to talk to their subjects.

There are a number of excellent texts on quantum mechanics at the undergraduate and graduate level. The following is a selection, beginning with the more elementary: A.P. FRENCH and EDWIN F. TAYLOR, *An Introduction to Quantum Physics* (1978); ALASTAIR I.M. RAE, *Quantum Mechanics*, 2nd ed. (1986); EUGEN MERZBACHER, *Quantum Mechanics*, 2nd ed. (1970); and ANTHONY SUDBERY, *Quantum Mechanics and the Particles of Nature: An Outline for Mathematicians* (1986), rather mathematical but including useful accounts and summaries of quantum metaphysics. RICHARD P. FEYNMAN, ROBERT B. LEIGHTON, and MATTHEW SANDS, *The Feynman Lectures on Physics*, vol. 3, *Quantum Mechanics* (1965), is a personal and stimulating look at the subject. A good introduction to quantum electrodynamics is RICHARD P. FEYNMAN, *QED: The Strange Theory of Light and Matter* (1985).

J.C. POLKINGHORNE, *The Quantum World* (1984); and JOHN GRIBBIN, *In Search of Schrödinger's Cat: Quantum Physics and Reality* (1984), are both highly readable and instructive books written at a popular level. BERNARD D'ESPAGNAT, *Conceptual Foundations of Quantum Mechanics*, 2nd ed. (1976), is a technical account of the fundamental conceptual problems involved. The proceedings of a conference, *New Techniques and Ideas in Quantum Measurement Theory*, ed. by DANIEL M. GREENBERGER (1986), contain a wide-ranging set of papers that deal with both the experimental and theoretical aspects of the measurement problem.

Applications are presented by H. HAKEN and H.C. WOLF, *Atomic and Quantum Physics: An Introduction to the Fundamentals of Experiment and Theory*, 2nd enlarged ed. (1987; originally published in German, 2nd rev. and enlarged ed., 1983); EMILIO SEGRÈ, *Nuclei and Particles: An Introduction to Nuclear and Subnuclear Physics*, 2nd rev. and enlarged ed. (1977, reissued 1980); DONALD H. PERKINS, *Introduction to High Energy Physics*, 3rd ed. (1987); CHARLES KITTEL, *Introduction to Solid State Physics*, 6th ed. (1986); and RODNEY LOUDON, *The Quantum Theory of Light*, 2nd ed. (1983). B.W. PETLEY, *The Fundamental Physical Constants and the Frontier of Measurement* (1985), gives a good account of present knowledge of the fundamental constants. (G.L.S.)

Medicine

Until the revolutionary scientific discoveries of the 19th and 20th centuries, medical practice was generally restricted to folk medicine and proscriptive religious and cultural tenets and limited by uneven knowledge and religious beliefs. As modern, empirical science took into its province the human body and its ills, "medicine" came to refer to the aggregate of scientific fields related to prevention and treatment of disease, as well as maintenance of health.

The many important discoveries in human anatomy and physiology, infectious and other diseases, drugs, and therapeutic procedures that took place during the 19th and 20th centuries have had a direct bearing on the important developments that occurred in the field of public health. Philanthropists, social activists, and government officials began to take a critical look at the living conditions of the poor, the ill, and the working class and joined in a positive effort toward reform and education. The laws enacted during this period gradually brought about clean water, safer living conditions, lower morbidities and mortalities, preventive medicine, and a decrease in the incidence of infectious diseases.

Medical education for physicians and other health-care professionals has generally kept pace with advances in medicine, especially in the developed countries. Each country in the world maintains its own requirements for medical degrees and licenses, and medical boards and councils have been formed to set standards and oversee the quality of medical education. Boards of certification have created stringent requirements that physicians must meet before they can practice a specialty, and they put great stress on the need for practicing physicians to continue their education as well.

A disparity in medical services is especially obvious be-

tween developed and developing countries. This disparity is the product of many problems, involving finances, communication, education, and cultural and religious views. International organizations such as the World Health Organization and the United Nations International Children's Fund make substantial progress in medical services by, among other things, considering the customs and culture in developing countries as an integral factor in the formation of the most efficient national health-care delivery system.

The concept of a general institution for the care of the sick can be traced to the Middle Ages, when monks set aside areas of their monasteries for this purpose. The modern general hospital fulfills a wide range of needs for the community both on an inpatient and an outpatient basis. The complexity and dynamic nature of medical practice is reflected in the growing number of specialty areas of the average general hospital. Teaching also has become a concern of many general hospitals, and, when designed as such, teaching hospitals provide services and knowledge often unavailable in other general hospitals.

Although religious and cultural laws as well as ethics have always been important in medicine, it was not until the 1960s that the modern legal system became intimately involved in medical practice. Advances in therapeutics and diagnostics have created legal and moral issues in such complex areas as abortion, euthanasia, and patients' rights. Physicians have had to reexamine the ethical tenets by which their predecessors practiced medicine and to develop new standards by which they may be judged.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, sections 423, 424, and 10/31, and the *Index*. (Ed.)

This article is divided into the following sections:

-
- | | | | |
|---|-----|--|-----|
| The history of medicine and surgery | 775 | China | |
| Medicine and surgery before 1800 | 775 | India | |
| Primitive medicine and folklore | | Other developing countries | |
| The ancient Middle East and Egypt | | Alternative or complementary medicine | |
| Traditional medicine and surgery in the Orient | | Special practices and fields of medicine | |
| The roots of Western medicine | | Specialties in medicine | |
| Christian and Muslim reservoirs of learning | | Teaching | |
| Medieval and Renaissance Europe | | Industrial medicine | |
| The Enlightenment | | Family health care | |
| The rise of scientific medicine in the 19th century | 782 | Geriatrics | |
| Physiology | | Public health practice | |
| Verification of the germ theory | | Military practice | |
| Discoveries in clinical medicine and anesthesia | | Clinical research | |
| Advances at the end of the century | | Historical notes | |
| Medicine in the 20th century | 783 | Clinical observation | |
| Infectious diseases and chemotherapy | | Drug research | |
| Immunology | | Surgery | |
| Endocrinology | | Screening procedures | |
| Vitamins | | Medical education | 801 |
| Malignant disease | | History of medical education | |
| Tropical medicine | | Modern patterns of medical education | |
| Surgery in the 20th century | 787 | Premedical education and admission to medical school | |
| The opening phase | | Undergraduate education | |
| World War I | | Postgraduate education | |
| Between the world wars | | Continuing education | |
| World War II and after | | Medical school faculty | |
| The practice of modern medicine | 791 | Requirements for practice | |
| Health care and its delivery | 791 | Economic aspects | |
| Organization of health services | | Scientific and international aspects | |
| Levels of health care | | Major medical institutions | 803 |
| Costs of health care | | Hospitals | 803 |
| Administration of primary health care | | History of hospitals | |
| Medical practice in developed countries | | The modern hospital | |
| Britain | | Comparison | |
| United States | | Ownership, control, and financing | |
| Russia | | The general hospital | |
| Japan | | Specialized health- and medical-care facilities | |
| Other developed countries | | Regional planning | |
| Medical practice in developing countries | | Public health services | 807 |

- History of public health
 - Beginnings in antiquity
 - The Middle Ages
 - The Renaissance
 - National developments in the 18th and 19th centuries
 - Developments from 1875
- Modern organizational and administrative patterns
- International organizations
 - Developed nations
 - Developing nations
- Progress in public health
 - Developed nations
 - Developing nations
- Clinics 813
 - Hospital clinics
 - Public health clinics
 - Private clinics
 - Health centres
 - Polyclinics
 - Family planning clinics
- Related fields 814
- Nursing 814
 - History of nursing
 - The practice of nursing
 - Kinds of nursing
 - Education
 - Licensing and registration
 - Organizations
 - Roles of international and multinational organizations
 - The International Red Cross
 - The World Health Organization
 - The European Economic Community
- Dentistry 817
 - History of dentistry
 - The practice of dentistry
 - Licensure requirements
 - Types of practice
 - Dental specialties and subspecialties
 - Dental education
 - Ancillary dental fields
 - Dental hygienists
 - Dental nurses and dental auxiliaries
 - Organizations
 - Pharmacy 821
 - History of pharmacy
 - The practice of pharmacy
 - Education
 - Licensing and regulation
 - Research
 - Organizations
 - Legal aspects of medicine 823
 - Maintenance of professional standards 823
 - History
 - Relationship of law and ethics
 - Law and medical practice 824
 - Government financing
 - License requirements for practice
 - The exclusive systems
 - The tolerant systems
 - The inclusive and integrated systems
 - Legal restrictions on practice
 - Determination of death
 - Termination of pregnancy
 - Sudden death
 - Public reporting
 - Legal redress
 - Patient's rights 826
 - Bibliography 827

THE HISTORY OF MEDICINE AND SURGERY

Medicine and surgery before 1800

PRIMITIVE MEDICINE AND FOLKLORE

Unwritten history is not easy to interpret, and, although much may be learned from a study of the drawings, bony remains, and surgical tools of early man, it is difficult to reconstruct his mental attitude toward the problems of disease and death. It seems probable that humans, as soon as they had reached the stage of reasoning, discovered, by the process of trial and error, which plants might be used as foods, which of them were poisonous, and which of them had some medicinal value. Folk medicine or domestic medicine, consisting largely in the use of vegetable products, or herbs, originated in this fashion and still persists.

But that is not the whole story. Man did not at first regard death and disease as natural phenomena. Common maladies, such as colds or constipation, were accepted as part of existence and dealt with by means of such herbal remedies as were available. Serious and disabling diseases, however, were placed in a very different category. These were of supernatural origin. They might be the result of a spell cast upon the victim by some enemy, visitation by a malevolent demon, or the work of an offended god who had either projected some object—a dart, a stone, a worm—into the body of the victim or had abstracted something, usually the soul of the patient. The treatment then applied was to lure the errant soul back to its proper habitat within the body or to extract the evil intruder, be it dart or demon, by counterspells, incantations, potions, suction, or other means.

One curious method of providing the disease with means of escape from the body was by making a hole, 2.5 to five centimetres across, in the skull of the victim—the practice of trepanning, or trephining. Trepanned skulls of prehistoric date have been found in Britain, France, and other parts of Europe and in Peru. Many of them show evidence of healing and, presumably, of the patient's survival. The practice still exists among primitive people in parts of Algeria, in Melanesia, and perhaps elsewhere, though it is fast becoming extinct.

Magic and religion played a large part in the medicine

of prehistoric or primitive man. Administration of a vegetable drug or remedy by mouth was accompanied by incantations, dancing, grimaces, and all the tricks of the magician. Therefore, the first doctors, or "medicine men," were witch doctors or sorcerers. The use of charms and talismans, still prevalent in modern times, is of ancient origin.

Apart from the treatment of wounds and broken bones, the folklore of medicine is probably the most ancient aspect of the art of healing, for primitive physicians showed their wisdom by treating the whole person, soul as well as body. Treatments and medicines that produced no physical effects on the body could nevertheless make a patient feel better when both medicine man and patient believed in their efficacy. This so-called placebo effect is applicable even in modern clinical medicine.

THE ANCIENT MIDDLE EAST AND EGYPT

The establishment of the calendar and the invention of writing marked the dawn of recorded history. The clues to early knowledge are scanty, consisting of clay tablets bearing cuneiform signs and seals that were used by physicians of ancient Mesopotamia. In the Louvre there is preserved a stone pillar on which is inscribed the Code of Hammurabi, who was a Babylonian king of the 18th century BC. This code includes laws relating to the practice of medicine, and the penalties for failure were severe. For example, "If the doctor, in opening an abscess, shall kill the patient, his hands shall be cut off"; if, however, the patient was a slave, the doctor was simply obliged to supply another slave.

The Greek historian Herodotus stated that every Babylonian was an amateur physician, since it was the custom to lay the sick in the street so that anyone passing by might offer advice. Divination, from the inspection of the liver of a sacrificed animal, was widely practiced to foretell the course of a disease. Little else is known regarding Babylonian medicine, and the name of not a single physician has survived.

When the medicine of ancient Egypt is examined, the picture becomes clearer. The first physician to emerge is Imhotep, chief minister to King Djoser in the 3rd millen-

Belief in
super-
natural
causes of
illness

The
physician
Imhotep

nium BC, who designed one of the earliest pyramids, the Step Pyramid at Saqqārah, and who was later regarded as the Egyptian god of medicine and identified with the Greek god Asclepius. Surer knowledge comes from the study of Egyptian papyri, especially the Ebers and Edwin Smith papyri discovered in the 19th century. The former is a list of remedies, with appropriate spells or incantations, while the latter is a surgical treatise on the treatment of wounds and other injuries.

Contrary to what might be expected, the widespread practice of embalming the dead body did not stimulate study of human anatomy. The preservation of mummies has, however, revealed some of the diseases suffered at that time, including arthritis, tuberculosis of the bone, gout, tooth decay, bladder stones, and gallstones; there is evidence too of the parasitic disease schistosomiasis, which remains a scourge still. There seems to have been no syphilis or rickets.

The search for information on ancient medicine leads naturally from the papyri of Egypt to Hebrew literature. Though the Bible contains little on the medical practices of Old Testament times, it is a mine of information on social and personal hygiene. The Jews were indeed pioneers in matters of public health. (D.J.G./P.Rh.)

TRADITIONAL MEDICINE AND SURGERY IN THE ORIENT

India. Indian medicine has a long history. Its earliest concepts are set out in the sacred writings called the Vedas, especially in the metrical passages of the Atharvaveda, which may possibly date as far back as the 2nd millennium BC. According to a later writer, the system of medicine called Āyurveda was received by a certain Dhanvantari from Brahma, and Dhanvantari was deified as the god of medicine. In later times his status was gradually reduced, until he was credited with having been an earthly king who died of snakebite.

The period of Vedic medicine lasted until about 800 BC. The Vedas are rich in magical practices for the treatment of diseases and in charms for the expulsion of the demons traditionally supposed to cause diseases. The chief conditions mentioned are fever (*takman*), cough, consumption, diarrhea, dropsy, abscesses, seizures, tumours, and skin diseases (including leprosy). The herbs recommended for treatment are numerous.

The golden age of Indian medicine, from 800 BC until about AD 1000, was marked especially by the production of the medical treatises known as the *Caraka-saṃhitā* and *Suśruta-saṃhitā*, attributed, respectively, to Caraka, a physician, and Suśruta, a surgeon. Estimates place the *Caraka-saṃhitā* in its present form as dating from the 1st century AD, although there were earlier versions. The *Suśruta-saṃhitā* probably originated in the last centuries BC and had become fixed in its present form by the 7th century AD. Of somewhat lesser importance are the treatises attributed to Vagbhata. All later writings on Indian medicine were based on these works.

Because Hindus were prohibited by their religion from cutting the dead body, their knowledge of anatomy was limited. The *Suśruta-saṃhitā* recommends that a body be placed in a basket and sunk in a river for seven days. On its removal the parts could be easily separated without cutting. As a result of these crude methods, the emphasis in Hindu anatomy was given first to the bones and then to the muscles, ligaments, and joints. The nerves, blood vessels, and internal organs were very imperfectly known.

The Hindus believed that the body contains three elementary substances, microcosmic representatives of the three divine universal forces, which they called spirit (air), phlegm, and bile (comparable to the humours of the Greeks). Health depends on the normal balance of these three elementary substances. The seven primary constituents of the body—blood, flesh, fat, bone, marrow, chyle, and semen—are produced by the action of the elementary substances. Semen was thought to be produced from all parts of the body and not from any individual part or organ.

Both Caraka and Suśruta state the existence of a large number of diseases (Suśruta says 1,120). Rough classifications of diseases are given. In all texts "fever," of which

numerous types are described, is regarded as important. Phthisis (wasting disease, especially pulmonary tuberculosis) was apparently prevalent, and the Hindu physicians knew the symptoms of cases likely to terminate fatally. Smallpox was common, and it is probable that smallpox inoculation was practiced.

Hindu physicians employed all five senses in diagnosis. Hearing was used to distinguish the nature of the breathing, alteration in voice, and the grinding sound produced by the rubbing together of broken ends of bones. They appear to have had a good clinical sense, and their discourses on prognosis contain acute references to symptoms that have grave import. Magical beliefs still persisted, however, until late in the classical period; thus, the prognosis could be affected by such fortuitous factors as the cleanliness of the messenger sent to fetch the physician, the nature of his conveyance, or the types of persons the physician met on his journey to the patient.

Dietetic treatment was important and preceded any medicinal treatment. Fats were much used, internally and externally. The most important methods of active treatment were referred to as the "five procedures": the administration of emetics, purgatives, water enemas, oil enemas, and sneezing powders. Inhalations were frequently administered, as were leeching, cupping, and bleeding.

The Indian materia medica was extensive and consisted mainly of vegetable drugs, all of which were from indigenous plants. Caraka knew 500 medicinal plants, and Suśruta knew 760. But animal remedies (such as the milk of various animals, bones, gallstones) and minerals (sulfur, arsenic, lead, copper sulfate, gold) were also employed. The physicians collected and prepared their own vegetable drugs. Among those that eventually appeared in Western pharmacopoeias were cardamom and cinnamon.

As a result of the strict religious beliefs of the Hindus, hygienic measures were important in treatment. Two meals a day were decreed, with indications of the nature of the diet, the amount of water to be drunk before and after the meal, and the use of condiments. Bathing and care of the skin were carefully prescribed, as were cleansing of the teeth with twigs from named trees, anointing of the body with oil, and the use of eyewashes.

In surgery, ancient Hindu medicine reached its zenith. Operations performed by Hindu surgeons included excision of tumours, incision and draining of abscesses, punctures to release fluid in the abdomen, extraction of foreign bodies, repair of anal fistulas, splinting of fractures, amputations, cesarean sections, and stitching of wounds.

A broad array of surgical instruments were used. According to Suśruta the surgeon should be equipped with 20 sharp and 101 blunt instruments of various descriptions. The instruments were largely of steel. Alcohol seems to have been used as a narcotic during operations, and bleeding was stopped by hot oils and tar.

In two types of operations especially, the Hindus were outstanding. Stone in the bladder (vesical calculus) was common in ancient India, and the surgeons frequently removed the stones by lateral lithotomy. They also introduced plastic surgery. Amputation of the nose was one of the prescribed punishments for adultery, and repair was carried out by cutting from the patient's cheek or forehead a piece of tissue of the required size and shape and applying it to the stump of the nose. The results appear to have been tolerably satisfactory, and the modern operation is certainly derived indirectly from this ancient source. Hindu surgeons also operated on cataracts by couching, or displacing the lens to improve vision.

China. The Chinese system of medicine is of great antiquity and is independent of any recorded external influences. According to legend, Emperor Huang Ti (the Yellow Emperor) wrote the canon of internal medicine called the *Nei ching* in the 3rd millennium BC; but there is some evidence that in its present form it dates from no earlier than the 3rd century BC. Most of the Chinese medical literature is founded on the *Nei ching*, and it is still regarded as a great authority. Other famous works are the *Mo ching* (known in the West as the "Pulse Classic"), composed about AD 300; and the *Golden Mirror*, a compilation, made about AD 1700, of medical writings of the

Hindu surgery

The Hindu view of the body

Han dynasty (202 BC–AD 220). European medicine began to obtain a footing in China early in the 19th century, but the native system is still widely practiced.

The
yin–yang
principle

Basic to traditional Chinese medicine is the dualistic cosmic theory of the yin and the yang. The yang, the male principle, is active and light and is represented by the heavens; the yin, the female principle, is passive and dark and is represented by the earth. The human body, like matter in general, is made up of five elements: wood, fire, earth, metal, and water. With these are associated other groups of five, such as the five planets, the five conditions of the atmosphere, the five colours, and the five tones. Health, character, and the success of all political and private ventures are determined by the preponderance, at the time, of the yin or the yang; and the great aim of ancient Chinese medicine is to control their proportions in the body.

The teachings of the religious sects forbade the mutilation of the dead human body; hence traditional anatomy rests on no sure scientific foundation. One of the most important writers on anatomy, Wang Ch'ing-jen, gained his knowledge from the inspection of dog-torn children who had died in a plague epidemic in AD 1798. Traditional Chinese anatomy is based on the cosmic system, which postulates the presence of such hypothetical structures as the 12 channels and the three so-called burning spaces. The body contains five organs (heart, lungs, liver, spleen, and kidneys), which store up but do not eliminate; and five viscera (such as the stomach, intestines, gallbladder, and bladder), which eliminate but do not store up. Each organ is associated with one of the planets, colours, tones, smells, and tastes. There are 365 bones and 365 joints in the body.

According to the physiology of traditional Chinese medicine, the blood vessels contain blood and air, in proportions varying with those of the yin and the yang. These two cosmic principles circulate in the 12 channels and control the blood vessels and hence the pulse. The *Nei ching* says that "the blood current flows continuously in a circle and never stops. It may be compared to a circle without beginning or end." On this insubstantial evidence it has been claimed that the Chinese anticipated Harvey's discovery of the circulation of the blood. Traditional Chinese pathology is also dependent on the theory of the yin and the yang; this led to an elaborate classification of diseases in which most of the types listed are without scientific foundation.

The pulse
doctrine

In diagnosis, detailed questions are asked about the history of the illness and about such things as the patient's taste, smell, and dreams. Conclusions are drawn from the quality of the voice, and note is made of the colour of the face and of the tongue. The most important part of the investigation, however, is the examination of the pulse. Wang Shu-ho, who wrote the "Pulse Classic," lived in the 3rd century BC, and innumerable commentaries were written on his work. The pulse is examined in several places, at different times, and with varying degrees of pressure. The operation may take as long as three hours. It is often the only examination made, and it is used both for diagnosis and for prognosis. Not only are the diseased organs ascertained but the time of death or recovery may be foretold.

The Chinese materia medica has always been extensive and consists of vegetable, animal (including human), and mineral remedies. There were famous herbals from ancient times; but all these, to the number of about 1,000, were embodied by Li Shih-chen in the compilation of *Pen-ts'ao kang-mu* (the "Great Pharmacopoeia") in the 16th century AD. This work, in 52 volumes, has been frequently revised and reprinted and is still authoritative. The use of drugs is mainly to restore the harmony of the yin and the yang and is also related to such matters as the five organs, the five planets, and the five colours. The art of prescribing is therefore complex.

Among the drugs taken over by Western medicine from the Chinese are rhubarb, iron (for anemia), castor oil, kaolin, aconite, camphor, and *Cannabis sativa* (Indian hemp). Chaulmoogra oil was used by the Chinese for leprosy from at least the 14th century, and about a century ago it began to be used for this purpose by Western physi-

cians. The herb mahuang (*Ephedra vulgaris*) has been used in China for at least 4,000 years, and the isolation of the alkaloid ephedrine from it has greatly improved the Western treatment of asthma and similar conditions.

The most famous and expensive of Chinese remedies is ginseng. Western analysis has shown that it has diuretic and other properties but is of doubtful value. In recent years reserpine, the active principle of the Chinese plant *Rauwolfia*, has been isolated; it is now effectively used in the treatment of high blood pressure and some emotional and mental conditions.

Hydrotherapy is probably of Chinese origin, since cold baths were used for fevers as early as 180 BC. The inoculation of smallpox matter, in order to produce a mild but immunizing attack of the disease, was practiced in China from ancient times and came to Europe about 1720. Another treatment is moxibustion, which consists in making a small, moistened cone (moxa) of powdered leaves of mugwort, or wormwood (*Artemisia* species), applying it to the skin, igniting it, and then crushing it into the blister so formed. Other substances are also used for the moxa. Dozens of these are sometimes applied at one sitting. The practice is often associated with acupuncture.

Acupuncture consists of the insertion into the skin and underlying tissues of a metal needle, either hot or cold. The theory is that the needle affects the distribution of the yin and the yang in the hypothetical channels and burning spaces of the body. The site of the insertion is chosen to affect a particular organ or organs. The practice of acupuncture dates from before 2500 BC and is peculiarly Chinese. Little of practical importance has been added since that date, although there have been many well-known treatises on the subject.

Acu-
puncture

A bronze model, c. AD 860, shows the hundreds of specified points for the insertion of the needle; this was the forerunner of countless later models and diagrams. The needles used are three to 24 centimetres (about one to nine inches) in length. They are often inserted with considerable force and after insertion may be agitated or screwed to the left or right. Acupuncture, often combined with moxibustion, is still widely used for many diseases, including fractures. Recently people in the Western world have turned to acupuncturists for relief from pain and other symptoms. There is some speculation that the treatment may trigger the brain to release morphinelike substances called endorphins, which presumably reduce the feeling of pain and its concomitant emotions.

Japan. The most interesting features of Japanese medicine are the extent to which it was derivative and the rapidity with which, after a slow start, it became westernized and scientific. In the early pre-Christian Era disease was regarded as sent by the gods or produced by the influence of evil spirits. Treatment and prevention were based largely on religious practices, such as prayers, incantations, and exorcism; at a later date drugs and bloodletting were also employed.

Beginning in AD 608, when young Japanese physicians were sent to China for a long period of study, Chinese influence on Japanese medicine was paramount. In 982, Tamba Yasuyori completed the 30-volume *Ishinhō*, the oldest Japanese medical work still extant. This work discusses diseases and their treatment, classified mainly according to the affected organs or parts. It is based entirely on older Chinese medical works, with the yin and yang concept underlying the theory of disease causation.

In 1570 a 15-volume medical work was published by Menase Dōsan, who also wrote at least five other works. In the most significant of these, the *Keitekishū* (a manual of the practice of medicine, 1574), diseases—or sometimes merely symptoms—are classified and described in 51 groups; the work is unusual in that it includes a section on the diseases of old age. Another distinguished physician and teacher of the period, Nagata Tokuhun, whose important books were the *I-no-ben* (1585) and the *Baika mujinzo* (1611), held that the chief aim of the medical art was to support the natural force, and consequently that it was useless to persist with stereotyped methods of treatment unless the physician had the cooperation of the patient.

European influences European medicine was introduced into Japan in the 16th century by Jesuit missionaries and again in the 17th century by Dutch physicians. Translations of European books on anatomy and internal medicine were made in the 18th century, and in 1836 an influential Japanese work on physiology appeared. In 1857 a group of Dutch-trained Japanese physicians founded a medical school in Edo (later Tokyo) that is regarded as the beginning of the medical faculty of the Imperial University of Tokyo.

During the last third of the 18th century it became government policy to westernize Japanese medicine, and great progress was made in the foundation of medical schools and the encouragement of research. Important medical breakthroughs by the Japanese followed, among them the discovery of the plague bacillus in 1894, the discovery of a dysentery bacillus in 1897, the isolation of adrenaline (epinephrine) in crystalline form in 1901, and the first experimental production of a tar-induced cancer in 1918.

(E.A.U./P.Rh.)

THE ROOTS OF WESTERN MEDICINE

Early Greece. The transition from magic to science was a gradual process that lasted for centuries, and there is little doubt that ancient Greece inherited much from Babylonia and Egypt, and even India and China. Twentieth-century readers of the Homeric tales the *Iliad* and the *Odyssey* may well be bewildered by the narrow distinction between gods and men among the characters and between historical fact and poetic fancy in the story. Two characters, the military surgeons Podaleirius and Machaon, are said to have been sons of Asclepius, the god of medicine. The divine Asclepius may have originated in a human Asclepius who lived about 1200 BC and is said to have performed many miracles of healing.

Asclepius was worshiped in hundreds of temples throughout Greece, the remains of which may still be seen at Epidaurus, Cos, Athens, and elsewhere. To these resorts, or hospitals, sick persons went for the healing ritual known as incubation, or temple sleep. They lay down to sleep in the dormitory, or *abaton*, and were visited in their dreams by Asclepius or by one of his priests, who gave advice. In the morning the patient often is said to have departed cured. There are at Epidaurus many inscriptions recording cures, though there is no mention of failures or deaths.

Diet, baths, and exercises played their part in the treatment, and it would appear that these temples were the prototype of modern health resorts. Situated in a peaceful spot, with gardens and fountains, each had its theatre for amusements and its stadium for athletic contests. The cult of incubation continued far into the Christian Era. In Greece, some of the Aegean islands, Sardinia, and Sicily, sick persons are still taken to spend a night in certain churches in the hope of a cure.

It was, however, the work of the early philosophers, rather than that of the priests of Asclepius, that impelled Greeks to refuse to be guided solely by supernatural influence and moved them to seek out for themselves the causes and reasons for the strange ways of nature. The 6th-century philosopher Pythagoras, whose chief discovery was the importance of numbers, also investigated the physics of sound, and his views influenced the medical thought of his time. In the 5th century BC Empedocles set forth the view that the universe is composed of four elements—fire, air, earth, and water; this conception led to the doctrine of the four bodily humours: blood; phlegm; choler, or yellow bile; and melancholy, or black bile. The maintenance of health was held to depend upon the harmony of the four humours.

Hippocrates. Medical thought had reached this stage and had partially discarded the conceptions based upon magic and religion by 460 BC, the year that Hippocrates is said to have been born. Although he has been called the father of medicine, little is known of his life, and there may, in fact, have been several men of this name; or Hippocrates may have been the author of only some, or none, of the books that make up the Hippocratic Collection (*Corpus Hippocraticum*). Ancient writers held that Hippocrates taught and practiced medicine in Cos, the island of his birth, and in other parts of Greece, including Athens, and that he died at an advanced age.

Whether Hippocrates was one man or several, the works attributed to him mark the stage in Western medicine where disease was coming to be regarded as a natural rather than a supernatural phenomenon and doctors were encouraged to look for physical causes of illness. Some of the works, notably the *Aphorismi* (*Aphorisms*), were used as textbooks until the 19th century. The first and best-known aphorism is, "Life is Short, Art long, Occasion sudden and dangerous, Experience deceitful, and Judgment difficult" (often shortened to the Latin tag, "Ars longa, vita brevis"). This is followed by brief comments on diseases and symptoms, many of which remain valid.

The thermometer and the stethoscope were not then known; nor, indeed, did Hippocrates employ any aid to diagnosis beyond his own powers of observation and logical reasoning. He had an extraordinary ability to foretell the course of a malady, and he laid more stress upon the expected outcome, or prognosis, of a disease than upon its identification, or diagnosis. He had no patience with the idea that disease was a punishment sent by the gods. Writing of epilepsy, then called "the sacred disease," he said, "It is not any more sacred than other diseases, but has a natural cause, and its supposed divine origin is due to man's inexperience. Every disease," he continued, "has its own nature, and arises from external causes."

Hippocrates noted the effect of food, of occupation, and especially of climate in causing disease, and one of his most interesting books, entitled *De aëre, aquis et locis* (*Air, Waters and Places*), would today be classed as a treatise on human ecology. Pursuing this line of thought, Hippocrates stated that "our natures are the physicians of our diseases" and advocated that this tendency to natural cure should be fostered. He laid much stress on diet and the use of few drugs. He knew well how to describe illness clearly and concisely and recorded failures as well as successes; he viewed disease with the eye of the naturalist and studied the entire patient in his environment.

Perhaps the greatest legacy of Hippocrates is the charter of medical conduct embodied in the so-called Hippocratic oath, which has been adopted as a pattern by physicians throughout the ages:

I swear by Apollo the physician, and Asclepius, and Health, and All-heal, and all the gods and goddesses . . . to reckon him who taught me this Art equally dear to me as my parents, to share my substance with him, and relieve his necessities if required; to look upon his offspring in the same footing as my own brothers, and to teach them this art, if they shall wish to learn it, without fee or stipulation; and that by precept, lecture, and every other mode of instruction, I will impart a knowledge of the Art to my own sons, and those of my teachers, and to disciples bound by a stipulation and oath according to the law of medicine, but to none others. I will follow that system of regimen which, according to my ability and judgment, I consider for the benefit of my patients, and abstain from whatever is deleterious and mischievous. I will give no deadly medicine to any one if asked, nor suggest any such counsel; and in like manner I will not give to a woman a pessary to produce abortion. . . . Into whatever houses I enter, I will go into them for the benefit of the sick, and will abstain from every voluntary act of mischief and corruption; and, further from the seduction of females or males, of freemen and slaves. Whatever, in connection with my professional practice or not, in connection with it, I see or hear, in the life of men, which ought not to be spoken of abroad, I will not divulge, as reckoning that all such should be kept secret.

Not strictly an oath, it was, rather, an ethical code or ideal, an appeal for right conduct. In one or other of its many versions, it has guided the practice of medicine throughout the world for more than 2,000 years.

Hellenistic and Roman medicine. In the following century the work of Aristotle, regarded as the first great biologist, was of inestimable value to medicine. A pupil of Plato at Athens and tutor to Alexander the Great, Aristotle studied the entire world of living things. He laid what can be identified as the foundations of comparative anatomy and embryology, and his views influenced scientific thinking for the next 2,000 years.

After the time of Aristotle, the centre of Greek culture shifted to Alexandria, where a famous medical school was

Cure by incubation

The Hippocratic oath

established in about 300 BC. There, the two best medical teachers were Herophilus, whose treatise on anatomy may have been the first of its kind, and Erasistratus, regarded by some as the founder of physiology. Erasistratus noted the difference between sensory and motor nerves but thought that the nerves were hollow tubes containing fluid and that air entered the lungs and heart and was carried through the body in the arteries. Alexandria continued as a centre of medical teaching even after the Roman Empire had attained supremacy over the Greek world, and medical knowledge remained predominantly Greek.

Asclepiades of Bithynia (born 124 BC) differed from Hippocrates in that he denied the healing power of nature and insisted that disease should be treated safely, speedily, and agreeably. An opponent of the humoral theory, he drew upon the atomic theory of the 5th-century Greek philosopher Democritus in advocating a doctrine of *strictum et laxum*—the attribution of disease to the contracted or relaxed condition of the solid particles that he believed make up the body. To restore harmony among the particles and thus effect cures, Asclepiades used typically Greek remedies: massage, poultices, occasional tonics, fresh air, and corrective diet. He gave particular attention to mental disease, clearly distinguishing hallucinations from delusions. He released the insane from confinement in dark cellars and prescribed a regimen of occupational therapy, soothing music, soporifics (especially wine), and exercises to improve the attention and memory.

Asclepiades did much to win acceptance for Greek medicine in Rome. Aulus Cornelius Celsus, the Roman nobleman who wrote *De medicina* about AD 30, gave a classic account of Greek medicine of the time, including descriptions of elaborate surgical operations. His book, overlooked in his day, enjoyed a wide reputation during the Renaissance.

During the early centuries of the Christian Era, Greek doctors thronged to Rome. The most illustrious of them was Galen, who began practicing there in AD 161. He acknowledged his debt to Hippocrates and followed the Hippocratic method, accepting the doctrine of the humours. He laid stress on the value of anatomy, and he virtually founded experimental physiology. Galen recognized that the arteries contain blood and not merely air. He showed how the heart sets the blood in motion in an ebb and flow fashion, but he had no idea that the blood circulates. Dissection of the human body was at that time illegal, so that he was forced to base his knowledge upon the examination of animals, particularly apes. A voluminous writer who stated his views forcibly and with confidence, he remained for centuries the undisputed authority from whom no one dared to differ.

Another influential physician of the 2nd century AD was Soranus of Ephesus, who wrote authoritatively on childbirth, infant care, and women's diseases. An opponent of abortion, he advocated numerous means of contraception. He also described how to assist a difficult delivery by turning the fetus in the uterus (podalic version), a life-saving technique that was subsequently lost sight of until it was revived in the 16th century.

Although the contribution of Rome to the practice of medicine was negligible compared with that of Greece, in matters of public health the Romans set the world a great example. The city of Rome had an unrivaled water supply. Gymnasiums and public baths were provided, and there was even domestic sanitation and adequate disposal of sewage. The army had its medical officers, public physicians were appointed to attend the poor, and hospitals were built; a Roman hospital excavated near Düsseldorf, Ger., was found to be strikingly modern in design.

CHRISTIAN AND MUSLIM RESERVOIRS OF LEARNING

After the fall of Rome, learning was no longer held in high esteem, experiment was discouraged, and originality became a dangerous asset. During the early Middle Ages medicine passed into the widely diverse hands of the Christian Church and Arab scholars.

Translators and saints. It is sometimes stated that the early Christian Church had an adverse effect upon medical progress. Disease was regarded as a punishment for

sin, and such chastening demanded only prayer and repentance. Moreover, the human body was held sacred and dissection was forbidden. But the infinite care and nursing bestowed upon the sick under Christian auspices must outweigh any intolerance shown toward medicine in the early days.

Perhaps the greatest service rendered to medicine by the church was the preservation and transcription of the classical Greek medical manuscripts. These were translated into Latin in many medieval monasteries, and the Nestorian Christians (an Eastern church) established a school of translators to render the Greek texts into Arabic. This famous school, and also a great hospital, were located at Jundi Shāhpūr in southwest Persia, where the chief physician was Jurjīs ibn Bukhtīshū, the first of a dynasty of translators and physicians that lasted for six generations. A later translator of great renown was Ḥunayn ibn Isḥāq, or Johannitus (born AD 809), whose translations were said to be worth their weight in gold.

About this time there appeared a number of saints whose names were associated with miraculous cures. Among the earliest of these were twin brothers, Cosmas and Damian, who suffered martyrdom (c. AD 303) and who became the patron saints of medicine. Other saints were invoked as powerful healers of certain diseases, such as St. Vitus for chorea (or St. Vitus' dance) and St. Anthony for erysipelas (or St. Anthony's fire). The cult of these saints was widespread in medieval times, and a later cult, that of St. Roch for plague, was widespread during the plague-ridden years of the 14th century.

Arabian medicine. A second reservoir of medical learning during those times was the great Muslim empire, which extended from Persia to Spain. Although it is customary to speak of Arabian medicine in describing this period, not all of the physicians were Arabs or natives of Arabia. Nor, indeed, were they all Muslims: some were Jews, some Christians, and they were drawn from all parts of the empire. One of the earliest figures was Rhazes, a Persian born in the last half of the 9th century near modern Tehrān, who wrote a voluminous treatise on medicine, *Kitāb al-hāwī* ("Comprehensive Book"), but whose most famous work, *De variolis et morbillis* (*A Treatise on the Smallpox and Measles*), distinguishes between these two diseases and gives a clear description of both.

Of later date was Avicenna (980–1037), also a Persian, who has been called the prince of physicians and whose tomb at Hamadan has become a place of pilgrimage. He could repeat the Qur'ān before he was 10 years old and at the age of 18 became court physician. His principal medical work, *al-Qānūn fī aṭ-ṭibb* (*The Canon of Medicine*), became a classic and was used at many medical schools—at Montpellier, Fr., as late as 1650—and reputedly is still used in the East.

The greatest contribution of Arabian medicine was in chemistry and in the knowledge and preparation of medicines. The chemists of that time were alchemists, and their pursuit was mainly a search for the philosopher's stone, which supposedly would turn common metals into gold. In the course of their experiments, however, numerous substances were named and characterized, and some were found to have medicinal value. Many drugs now in use are of Arab origin, as are such processes as distillation and sublimation.

At that period, and indeed throughout most historical times, surgery was considered inferior to medicine, and surgeons were held in low regard. The renowned Spanish surgeon Abū al-Qāsim (Albucasis), however, did much to raise the status of surgery in Córdoba, an important centre of commerce and culture with a hospital and medical school equal to those of Cairo and Baghdad. A careful and conservative practitioner, he wrote the first illustrated surgical text, which held wide influence in Europe for centuries.

Another great doctor of Córdoba, born in the 12th century, just as the sun of Arabian culture was setting, was the Jewish philosopher Maimonides. Banished from the city because he would not become a Muslim, he eventually went to Cairo, where the law was more lenient and where he acquired a reputation so high that he became physician

Preservation of Greek medical manuscripts

Avicenna

Galen

to Saladin, the Saracen leader. (He was the original of El Hakim in Sir Walter Scott's *Talisman*.) A few of his works, written in Arabic, were eventually translated into Latin and printed; perhaps the best known is *Yad-Hachazakah*, or *The Code of Maimonides*, an ethical guide still esteemed in Jewish medical circles.

MEDIEVAL AND RENAISSANCE EUROPE

Salerno and the medical schools. At about the same time that Arabian medicine flourished, the first organized medical school in Europe was established at Salerno, in southern Italy. Although the school of Salerno produced no brilliant genius and no startling discovery, it was the outstanding medical institution of its time and the parent of the great medieval schools soon to be founded at Montpellier and Paris, in France, and at Bologna and Padua, in Italy. Salerno drew scholars from near and far. Remarkably liberal in some of its views, Salerno admitted women as medical students. The school owed much to the enlightened Holy Roman emperor Frederick II, who decreed in 1221 that no one should practice medicine until he had been publicly approved by the masters of Salerno.

The Salernitan school also produced a literature of its own; the best-known work, of uncertain date and of composite authorship, was the *Regimen Sanitatis Salernitanum* ("Salernitan Guide to Health"). Written in verse, it has appeared in numerous editions and has been translated into many languages. Among its oft-quoted couplets is the following:

Use three physicians still, first Doctor Quiet,
Next Doctor Merryman, and Doctor Diet.

Salerno yielded its place as the premier medical school of Europe to Montpellier in about 1200. John of Gaddesden, the model for the "doctour of physick" in Chaucer's *Canterbury Tales*, was one of the English students there. That he relied upon astrology and upon the doctrine of the humours is evident from Chaucer's description:

Well could he guess the ascending of the star
Wherein his patient's fortunes settled were,
He knew the course of every malady,
Were it of cold or heat or moist or dry.

Medieval physicians analyzed symptoms, examined excreta, and made their diagnoses. Then they might prescribe diet, rest, sleep, exercise, or baths; or they could administer emetics and purgatives or bleed the patient. Surgeons could treat fractures and dislocations, repair hernias, and perform amputations and a few other operations. Some of them prescribed opium, mandragora, or alcohol to deaden pain. Childbirth was left to midwives, who relied on folklore and tradition.

Great hospitals were established during the Middle Ages by religious foundations, and infirmaries were attached to abbeys, monasteries, priories, and convents. Doctors and nurses in these institutions were members of religious orders and combined spiritual with physical healing.

The spread of new learning. Among the teachers of medicine in the medieval universities there were many who clung to the past, but there were not a few who determined to explore new lines of thought. The new learning of the Renaissance, born in Italy, grew and expanded slowly. Two great 13th-century scholars who influenced medicine were Roger Bacon, an active observer and tireless experimenter, and Albertus Magnus, a distinguished philosopher and scientific writer.

About this time Mondino dei Liucci taught at Bologna. Prohibitions against human dissection were slowly lifting, and Mondino performed his own dissections rather than following the customary procedure of entrusting the task to a menial. Although he perpetuated the errors of Galen, his *Anothomia*, published in 1316, was the first practical manual of anatomy. Foremost among the surgeons of the day was Guy de Chauliac, a physician to three popes at Avignon. His *Chirurgia magna* ("Great Surgery"), based on observation and experience, had a profound influence upon the progress of surgery.

The Renaissance in the 14th, 15th, and 16th centuries was much more than just a reviving of interest in Greek and Roman culture; it was rather a change of outlook, an

eagerness for discovery, a desire to escape from the limitations of tradition and to explore new fields of thought and action. In medicine, it was perhaps natural that anatomy and physiology, the knowledge of the human body and its workings, should be the first aspects of medical learning to receive attention from those who realized the need for reform.

It was in 1543 that Andreas Vesalius, a young Belgian professor of anatomy at the University of Padua, published *De humani corporis fabrica* ("On the Structure of the Human Body"). Based on his own dissections, this seminal work corrected many of Galen's errors. By his scientific observations and methods, Vesalius showed that Galen could no longer be regarded as the final authority. His work at Padua was continued by Gabriel Fallopius and, later, by Hieronymus Fabricius ab Aquapendente; it was his work on the valves in the veins, *De venarum ostioliis* (1603), that suggested to his pupil William Harvey his revolutionary theory of the circulation of the blood, one of the great medical discoveries.

Surgery profited from the new outlook in anatomy, and the great reformer Ambroise Paré dominated the field in the 16th century. Paré was surgeon to four kings of France, and he has deservedly been called the father of modern surgery. In his autobiography, written after he had retired from 30 years of service as an army surgeon, Paré described how he had abolished the painful practice of cautery to stop bleeding and used ligatures and dressings instead. His favourite expression, "I dressed him; God healed him," is characteristic of this humane and careful doctor.

In Britain during this period surgery, which was performed by barber-surgeons, was becoming regulated and organized under royal charters. Companies were thus formed that eventually became the royal colleges of surgeons in Scotland and England. Physicians and surgeons united in a joint organization in Glasgow, and a college of physicians was founded in London.

The 16th-century medical scene was enlivened by the enigmatic physician and alchemist who called himself Paracelsus. Born in Switzerland, he traveled extensively throughout Europe, gaining medical skills and practicing and teaching as he went. In the tradition of Hippocrates, Paracelsus stressed the power of nature to heal; but unlike Hippocrates he believed also in the power of supernatural forces, and he violently attacked the medical treatments of his day. Eager for reform, he allowed his intolerance to outweigh his discretion, as when he prefaced his lectures at Basel by publicly burning the works of Avicenna and Galen. The authorities and medical men were understandably outraged. Widely famous in his time, Paracelsus remains a controversial figure to this day. Despite his turbulent career, however, he did attempt to bring a more rational approach to diagnosis and treatment, and he introduced the use of chemical drugs in place of herbal remedies.

A contemporary of Paracelsus, Girolamo Fracastoro of Italy was a scholar cast from a very different mold. His account of the disease syphilis, entitled *Syphilis sive morbus Gallicus* (1530; "Syphilis or the French Disease"), was written in verse. Although Fracastoro called syphilis the French disease, others called it the Neapolitan disease, for it was said to have been brought to Naples from America by the sailors of Christopher Columbus. Its origin is still questioned, however. Fracastoro was interested in epidemic infection, and he offered the first scientific explanation of disease transmission. In his great work, *De contagione et contagiosis morbis* (1546), he theorized that the seeds of certain diseases are imperceptible particles transmitted by air or by contact.

THE ENLIGHTENMENT

In the 17th century the natural sciences moved forward on a broad front. There were attempts to grapple with the nature of science, as expressed in the works of thinkers like Francis Bacon, Descartes, and Newton. New knowledge of chemistry superseded the theory that all things are made up of earth, air, fire, and water, and the old Aristotelian ideas began to be discarded. The supreme 17th-century

Vesalius

Mont-
pellier

Fracastoro
on syphilis

achievement in medicine was Harvey's explanation of the circulation of blood.

Harvey and the experimental method. Born in Folkestone, Eng., William Harvey studied at Cambridge University and then spent several years at Padua, where he came under the influence of Fabricius. He established a successful medical practice in London and by precise observation and scrupulous reasoning developed his theory of circulation. In 1628 he published his classic book *Exercitatio Anatomica de Motu Cordis et Sanguinis in Animalibus* (*Concerning the Motion of the Heart and Blood*), often called *De Motu Cordis*.

That the book aroused controversy is not surprising. There were still many who adhered to the teaching of Galen that the blood follows an ebb and flow movement in the blood vessels. Harvey's work was the result of many careful experiments, but few of his critics took the trouble to repeat the experiments, simply arguing in favour of the older view. His second great book, *Exercitationes de generatione animalium* ("Experiments Concerning Animal Generation"), published in 1651, laid the foundation of modern embryology.

Harvey's discovery of the circulation of the blood was a landmark of medical progress; the new experimental method by which the results were secured was as noteworthy as the work itself. Following the method described by the philosopher Francis Bacon, he drew the truth from experience and not from authority.

There was one gap in Harvey's argument: he was obliged to assume the existence of the capillary vessels that conveyed the blood from the arteries to the veins. This link in the chain of evidence was supplied by Marcello Malpighi of Bologna (who was born in 1628, the year of publication of *De Motu Cordis*). With a primitive microscope Malpighi saw a network of tiny blood vessels in the lung of a frog. Harvey also failed to show why the blood circulated. After Robert Boyle had shown that air is essential to animal life, it was Richard Lower who traced the interaction between air and the blood. Eventually the importance of oxygen, which was confused for a time by some as phlogiston, was revealed, although it was not until the late 18th century that the great chemist Antoine-Laurent Lavoisier discovered the essential nature of oxygen and clarified its relation to respiration.

The microscope

Although the compound microscope had been invented slightly earlier, probably in Holland, its development, like that of the telescope, was the work of Galileo. He was the first to insist upon the value of measurement in science and in medicine, thus replacing theory and guesswork with accuracy. The great Dutch microscopist Antonie van Leeuwenhoek devoted his long life to microscopical studies and was probably the first to see and describe bacteria, reporting his results to the Royal Society of London. In England, Robert Hooke, who was Boyle's assistant and curator to the Royal Society, published his *Micrographia* in 1665, which discussed and illustrated the microscopic structure of a variety of materials.

The futile search for an easy system. Several attempts were made in the 17th century to discover an easy system that would guide the practice of medicine. A substratum of superstition still remained. Richard Wiseman, surgeon to Charles II, affirmed his belief in the "royal touch" as a cure for king's evil, or scrofula, while even the learned English physician Thomas Browne stated that witches really existed. There was, however, a general desire to discard the past and adopt new ideas.

The view of the French philosopher René Descartes that the human body is a machine and that it functions mechanically had its repercussions in medical thought. One group adopting this explanation called themselves the iatrophysicists; another school, preferring to view life as a series of chemical processes, were called iatrochemists. Santorio Santorio, working at Padua, was an early exponent of the iatrophysical view and a pioneer investigator of metabolism. He was especially concerned with the measurement of what he called "insensible perspiration," described in his book *De statica medicina* (1614; "On Medical Measurement"). Another Italian, who developed the idea still further, was Giovanni Alfonso Borelli, a

professor of mathematics at Pisa University, who gave his attention to the mechanics and statics of the body and to the physical laws that govern its movements.

The iatrochemical school was founded at Brussels by Jan Baptist van Helmont, whose writings are tinged with the mysticism of the alchemist. A more logical and intelligible view of iatrochemistry was advanced by Franciscus Sylvius, at Leiden; and in England a leading exponent of the same school was Thomas Willis, who is better known for his description of the brain in his *Cerebri anatome nervorumque descriptio et usus* ("Anatomy of the Brain and Descriptions and Functions of the Nerves"), published in 1664 and illustrated by Christopher Wren.

It soon became apparent that no easy road to medical knowledge and practice was to be found along these channels and that the best method was the age-old system of straightforward clinical observation initiated by Hippocrates. The need for a return to these views was strongly urged by Thomas Sydenham, well named "the English Hippocrates." Sydenham was not a voluminous writer and, indeed, had little patience with book learning in medicine; nevertheless he gave excellent descriptions of the phenomena of disease. His greatest service, much needed at the time, was to divert physicians' minds from speculation and lead them back to the bedside, where the true art of medicine could be studied.

Medicine in the 18th century. Even in the 18th century the search for a simple way of healing the sick continued. In Edinburgh the writer and lecturer John Brown expounded his view that there were only two diseases, sthenic (strong) and asthenic (weak), and two treatments, stimulant and sedative; his chief remedies were alcohol and opium. Lively and heated debates took place between his followers, the Brunonians, and the more orthodox Cullenians (followers of William Cullen, a professor of medicine at Glasgow), and the controversy spread to the medical centres of Europe.

At the opposite end of the scale, at least in regard to dosage, was Samuel Hahnemann, of Leipzig, the originator of homeopathy, a system of treatment involving the administration of minute doses of drugs whose effects resemble the effects of the disease being treated. His ideas had a salutary effect upon medical thought at a time when prescriptions were lengthy and doses were large, and his system has had many followers.

By the 18th century the medical school at Leiden had grown to rival that of Padua, and many students were attracted there from abroad. Among them was John Monro, an army surgeon, who resolved that his native city of Edinburgh should have a similar medical school. He specially educated his son Alexander with a view to having him appointed professor of anatomy, and the bold plan was successful. Alexander Monro studied at Leiden under Hermann Boerhaave, the central figure of European medicine and the greatest clinical teacher of his time. Subsequently, three generations of Alexander Monros taught anatomy at Edinburgh University over a continuous period of 126 years. Medical education was increasingly incorporated into the universities of Europe, and Edinburgh became the leading academic centre for medicine in Britain.

Leiden and Edinburgh

In 18th-century London, Scottish doctors were the leaders in surgery and obstetrics. The noted teacher John Hunter conducted extensive researches in comparative anatomy and physiology, founded surgical pathology, and raised surgery to the level of a respectable branch of science. His brother William Hunter, an eminent teacher of anatomy, became famous as an obstetrician. Male doctors were now attending women in childbirth, and the leading obstetrician in London was William Smellie. His well-known *Treatise on the Theory and Practice of Midwifery*, published in three volumes in 1752-64, contained the first systematic discussion on the safe use of obstetrical forceps, which have since saved countless lives. Smellie placed midwifery on a sound scientific footing and helped to establish obstetrics as a recognized medical discipline.

The science of modern pathology also had its beginnings in this century. Giovanni Battista Morgagni, of Padua, in 1761 published his massive work *De sedibus et causis morborum* (*The Seats and Causes of Diseases Investigated*

by *Anatomy*), a description of the appearances found by postmortem examination of almost 700 cases, in which he attempted to correlate the findings after death with the clinical picture in life.

On the basis of work begun in the 18th century, René Laënnec, a native of Brittany, who practiced medicine in Paris, invented a simple stethoscope, or *cylindre*, as it was originally called. In 1819 he wrote a treatise, *De l'auscultation médiate* ("On Mediate Auscultation"), describing many of the curious sounds in the heart and lungs that are revealed by the instrument. Meanwhile a Viennese physician, Leopold Auenbrugger, discovered another method of investigating diseases of the chest, that of percussion. The son of an innkeeper, he is said to have conceived the idea of tapping with the fingers when he recalled that he had used this method to gauge the level of the fluid contents of his father's casks.

One highly significant medical advance, late in the century, was vaccination. Smallpox, disfiguring and often fatal, was widely prevalent. Inoculation, which had been practiced in the East, was popularized in England in 1721–22 by Lady Mary Wortley Montagu, who is best known for her letters. She observed the practice in Turkey, where it produced a mild form of the disease, thus securing immunity, although not without danger. The next step was taken by Edward Jenner, a country practitioner who had been a pupil of John Hunter. In 1796 Jenner began inoculations with material from cowpox (the bovine form of the disease); and when he later inoculated the same subject with smallpox, the disease did not appear. This procedure—vaccination—has been responsible for eradicating the disease.

Smallpox
vaccination

Public health and hygiene were receiving more attention during the 18th century. Population statistics began to be kept, and suggestions arose concerning health legislation. Hospitals were established for a variety of purposes. In Paris, Philippe Pinel initiated bold reforms in the care of the mentally ill, releasing them from their chains and discarding the long-held notion that insanity was caused by demon possession.

Conditions improved for sailors and soldiers as well. James Lind, a British naval surgeon from Edinburgh, recommended fresh fruits and citrus juices to prevent scurvy, a remedy discovered by the Dutch in the 16th century. When the British navy adopted Lind's advice—decades later—this deficiency disease was eliminated. In 1752 a Scotsman, John Pringle, published his classic *Observations on the Diseases of the Army*, which contained numerous recommendations for the health and comfort of the troops. Serving with the British forces during the War of the Austrian Succession, he suggested in 1743 that military hospitals on both sides should be regarded as sanctuaries; this plan eventually led to the establishment of the Red Cross organization in 1864.

Two pseudoscientific doctrines relating to medicine emerged from Vienna in the latter part of the century and attained wide notoriety. Mesmerism, a belief in "animal magnetism" sponsored by Franz Anton Mesmer, probably owed any therapeutic value it had to suggestions given while the patient was under hypnosis. Phrenology, propounded by Franz Joseph Gall, held that the contours of the skull were a guide to an individual's mental faculties and character traits; this theory remained popular throughout the 19th century.

At the same time, sound scientific thinking was making steady progress, and advances in physics, chemistry, and the biological sciences were converging to form a rational scientific basis for every branch of clinical medicine. New knowledge disseminated throughout Europe and traveled across the sea, where centres of medical excellence were being established in America.

The rise of scientific medicine in the 19th century

The portrayal of the history of medicine becomes more difficult in the 19th century. Discoveries multiply, and the number of eminent doctors is so great that the history is apt to become a series of biographies. Nevertheless, it

is possible to discern the leading trends in modern medical thought.

PHYSIOLOGY

By the beginning of the 19th century, the structure of the human body was almost fully known, due to new methods of microscopy and of injections. Even the body's microscopic structure was understood. But as important as anatomical knowledge was an understanding of physiological processes, which were rapidly being elucidated, especially in Germany. There, physiology became established as a distinct science under the guidance of Johannes Müller, who was a professor at Bonn and then at the University of Berlin. An energetic worker and an inspiring teacher, he described his discoveries in a famous textbook, *Handbuch der Physiologie des Menschen* ("Manual of Human Physiology"), published in the 1830s.

Among Müller's illustrious pupils were Hermann von Helmholtz, who made significant discoveries relating to sight and hearing and who invented the ophthalmoscope; and Rudolf Virchow, one of the century's great medical scientists, whose outstanding achievement was his conception of the cell as the centre of all pathological changes. Virchow's work *Die Cellularpathologie*, published in 1858, gave the deathblow to the outmoded view that disease is due to an imbalance of the four humours.

In France the most brilliant physiologist of the time was Claude Bernard, whose many important discoveries were the outcome of carefully planned experiments. His researches clarified the role of the pancreas in digestion, revealed the presence of glycogen in the liver, and explained how the contraction and expansion of the blood vessels are controlled by vasomotor nerves. He proposed the concept of the internal environment—the chemical balance in and around the cells—and the importance of its stability. His *Introduction à l'étude de la médecine expérimentale* (1865; *An Introduction to the Study of Experimental Medicine*) is still worthy of study by all who undertake research.

VERIFICATION OF THE GERM THEORY

Perhaps the overarching medical advance of the 19th century, certainly the most spectacular, was the conclusive demonstration that certain diseases, as well as the infection of surgical wounds, were directly caused by minute living organisms. This discovery changed the whole face of pathology and effected a complete revolution in the practice of surgery.

The idea that disease was caused by entry into the body of imperceptible particles was of ancient date. It had been expressed by the Roman encyclopaedist Varro as early as 100 BC, by Fracastoro in 1546, by Athanasius Kircher and Pierre Borel about a century later, and by Francesco Redi, who in 1684 wrote his *Osservazioni intorno agli animali viventi che si trovano negli animali viventi* ("Observations on Living Animals Which Are to Be Found Within Other Living Animals"), in which he sought to disprove the idea of spontaneous generation. Everything must have a parent, he wrote; only life produces life. A 19th-century pioneer in this field, regarded by some as founder of the parasitic theory of infection, was Agostino Bassi of Italy, who showed that a disease of silkworms was caused by a fungus that could be destroyed by chemical agents.

The main credit for establishing the science of bacteriology must be accorded to the French chemist Louis Pasteur. It was Pasteur who, by a brilliant series of experiments, proved that the fermentation of wine and the souring of milk are caused by living microorganisms. His work led to the pasteurization of milk and solved problems of agriculture and industry as well as those of animal and human diseases. He successfully employed inoculations to prevent anthrax in sheep and cattle, chicken cholera in fowl, and finally rabies in humans and dogs. The latter resulted in the widespread establishment of Pasteur institutes.

From Pasteur, Joseph Lister derived the concepts that enabled him to introduce the antiseptic principle into surgery. In 1865 Lister, a professor of surgery at Glasgow University, began placing an antiseptic barrier of carbolic acid between the wound and the germ-containing atmosphere. Infections and deaths fell dramatically, and his

The work of Johannes Müller and his students

Influence of Pasteur

pioneering work led to more refined techniques of sterilizing the surgical environment.

Obstetrics had already been robbed of some of its terrors by Alexander Gordon at Aberdeen, Scot., Oliver Wendell Holmes at Boston, and Ignaz Semmelweis at Vienna and Pest (Budapest), who advocated disinfection of the hands and clothing of midwives and medical students who attended confinements. These measures produced a marked reduction in cases of puerperal fever, the bacterial scourge of women following childbirth.

Another pioneer in bacteriology was the German physician Robert Koch, who showed how bacteria could be cultivated, isolated, and examined in the laboratory. A meticulous investigator, Koch discovered the organisms of tuberculosis, in 1882, and cholera, in 1883. By the end of the century many other disease-producing microorganisms had been identified.

DISCOVERIES IN CLINICAL MEDICINE AND ANESTHESIA

There was perhaps some danger that in the search for bacteria other causes of disease would escape detection. Many physicians, however, were working along different lines in the 19th century. Among them were a group attached to Guy's Hospital, in London: Richard Bright, Thomas Addison, and Sir William Gull. Bright contributed significantly to the knowledge of kidney diseases, including Bright's disease, and Addison gave his name to disorders of the adrenal glands and the blood. Gull, a famous clinical teacher, left a legacy of pithy aphorisms that might well rank with those of Hippocrates.

In Dublin Robert Graves and William Stokes introduced new methods in clinical diagnosis and medical training; while in Paris a leading clinician, Pierre-Charles-Alexandre Louis, was attracting many students from America by the excellence of his teaching. By the early 19th century the United States was ready to send back the results of its own researches and breakthroughs. In 1809, in a small Kentucky town, Ephraim McDowell boldly operated on a woman—without anesthesia or antisepsis—and successfully removed a large ovarian tumour. William Beaumont, in treating a shotgun wound of the stomach, was led to make many original observations that were published in 1833 as *Experiments and Observations on the Gastric Juice and the Physiology of Digestion*.

The most famous contribution by the United States to medical progress at this period was undoubtedly the introduction of general anesthesia, a procedure that not only liberated the patient from the fearful pain of surgery but also enabled the surgeon to perform more extensive operations. The discovery was marred by controversy. Crawford Long, Gardner Colton, Horace Wells, and Charles Jackson are all claimants for priority; some used nitrous oxide gas, and others employed ether, which was less capricious. There is little doubt, however, that it was William Thomas Morton who, on Oct. 16, 1846, at Massachusetts General Hospital, in Boston, first demonstrated before a gathering of physicians the use of ether as a general anesthetic. The news quickly reached Europe, and general anesthesia soon became prevalent in surgery. At Edinburgh, the professor of midwifery, James Young Simpson, had been experimenting upon himself and his assistants, inhaling various vapours with the object of discovering an effective anesthetic. In November 1847 chloroform was tried with complete success, and soon it was preferred to ether and became the anesthetic of choice.

ADVANCES AT THE END OF THE CENTURY

While antisepsis and anesthesia placed surgery on an entirely new footing, similarly important work was carried out in other fields of study, such as parasitology and disease transmission. Patrick Manson, a British pioneer in tropical medicine, showed in China, in 1877, how insects can carry disease and how the embryos of the *Filaria* worm, which can cause elephantiasis, are transmitted by the mosquito. Manson explained his views to a British army surgeon, Ronald Ross, then working on the problem of malaria, and Ross discovered the malarial parasite in the stomach of the *Anopheles* mosquito in 1897.

In Cuba, Carlos Finlay expressed the view, in 1881, that

yellow fever is carried by the *Stegomyia* mosquito. Following his lead, the Americans Walter Reed, William Gorgas, and others were able to conquer the scourge of yellow fever in Panama and made possible the completion of the Panama Canal by reducing the death rate there from 176 per 1,000 to 6 per 1,000.

Other victories in preventive medicine ensued, because the maintenance of health was now becoming as important a concern as the cure of disease; and the 20th century was to witness the evolution and progress of national health services in a number of countries. In addition, spectacular advances in diagnosis and treatment followed the discovery of X rays by Wilhelm Conrad Röntgen, in 1895, and of radium by Pierre and Marie Curie in 1898. Before the turn of the century, too, the vast new field of psychiatry had been opened up by Sigmund Freud. The tremendous increase in scientific knowledge during the 19th century radically altered and expanded the practice of medicine. Concern for upholding the quality of services led to the establishment of public and professional bodies to govern the standards for medical training and practice.

(D.J.G./P.Rh.)

Medicine in the 20th century

The 20th century has produced such a plethora of discoveries and advances that in some ways the face of medicine has changed out of all recognition. In 1901, for instance, in the United Kingdom the expectation of life at birth, a primary indicator of the effect of health care on mortality (but also reflecting the state of health education, housing, and nutrition), was 48 years for males and 51.6 years for females. After steady increases, by the 1980s life expectancy had reached 71.4 years for males and 77.2 years for females. Other industrialized nations showed similar dramatic increases. Indeed, the outlook has so altered that, with the exception of diseases such as cancer and AIDS, attention has become focused on morbidity rather than mortality, and the emphasis has changed from keeping people alive to keeping them fit.

The rapid progress of medicine in this era was reinforced by enormous improvements in communication between scientists throughout the world. Through publications, conferences, and—later—computers and electronic media, they freely exchanged ideas and reported on their endeavours. No longer was it common for an individual to work in isolation. Although specialization increased, teamwork became the norm. It consequently has become more difficult to ascribe medical accomplishments to particular individuals.

In the first half of the century, emphasis continued to be placed on combating infection, and notable landmarks were also attained in endocrinology, nutrition, and other areas. In the years following World War II, insights derived from cell biology altered basic concepts of the disease process; new discoveries in biochemistry and physiology opened the way for more precise diagnostic tests and more effective therapies; and spectacular advances in biomedical engineering enabled the physician and surgeon to probe into the structures and functions of the body by noninvasive imaging techniques like ultrasound (sonar), computerized axial tomography (CAT), and nuclear magnetic resonance (NMR). With each new scientific development, medical practices of just a few years earlier became obsolete.

INFECTIOUS DISEASES AND CHEMOTHERAPY

In the years following the turn of the century, ongoing research concentrated on the nature of infectious diseases and their means of transmission. Increasing numbers of pathogenic organisms were discovered and classified. Some, such as the rickettsias, which cause diseases like typhus, were smaller than bacteria; some were larger, such as the protozoans that engender malaria and other tropical diseases. The smallest to be identified were the viruses, producers of many discases, among them mumps, measles, German measles, and poliomyelitis; and in 1910 Peyton Rous showed that a virus could also cause a malignant tumour, a sarcoma in chickens.

Varieties of micro-organisms

The 19th century

Insect transmission of disease

There was still little to be done for the victims of most infectious organisms beyond drainage, poultices, and ointments, in the case of local infections, and rest and nourishment for severe diseases. The search for treatments aimed at both vaccines and chemical remedies.

Ehrlich and arsphenamine. Germany was well to the forefront in medical progress. The scientific approach to medicine had been developed there long before it spread to other countries, and postgraduates flocked to German medical schools from all over the world. The opening decade of the 20th century has been well described as the golden age of German medicine. Outstanding among its leaders was Paul Ehrlich.

While still a student, Ehrlich carried out some work on lead poisoning from which he evolved the theory that was to guide much of his subsequent work—that certain tissues have a selective affinity for certain chemicals. He experimented with the effects of various chemical substances on disease organisms. In 1910, with his colleague Sahachiro Hata, he conducted tests on arsphenamine, once sold under the commercial name Salvarsan. Their success inaugurated the chemotherapeutic era, which was to revolutionize the treatment and control of infectious diseases. Salvarsan, a synthetic preparation containing arsenic, is lethal to the microorganism responsible for syphilis. Until the introduction of penicillin, Salvarsan or one of its modifications remained the standard treatment of syphilis and went far toward bringing this social and medical scourge under control.

Sulfonamide drugs. In 1932 the German bacteriologist Gerhard Domagk announced that the red dye Prontosil is active against streptococcal infections in mice and humans. Soon afterward French workers showed that its active antibacterial agent is sulfanilamide. In 1936 the English physician Leonard Colebrook and his colleagues provided overwhelming evidence of the efficacy of both Prontosil and sulfanilamide in streptococcal septicemia (bloodstream infection), thereby ushering in the sulfonamide era. New sulfonamides, which appeared with astonishing rapidity, had greater potency, wider antibacterial range, or lower toxicity. Some stood the test of time; others, like the original sulfanilamide and its immediate successor, sulfapyridine, were replaced by safer and more powerful successors.

Antibiotics. *Penicillin.* A dramatic episode in medical history occurred in 1928, when Alexander Fleming noticed the inhibitory action of a stray mold on a plate culture of staphylococcus bacteria in his laboratory at St. Mary's Hospital, London. Many other bacteriologists must have made the observation, but none had realized the possible implications. The mold was a strain of *Penicillium*—*P. notatum*—which gave its name to the now-famous drug penicillin. In spite of his conviction that penicillin was a potent antibacterial agent, Fleming was unable to carry his work to fruition, mainly because biochemists at the time were unable to isolate it in sufficient quantities or in a sufficiently pure form to allow its use on patients.

Ten years later Howard Florey, Ernst Chain, and their colleagues at Oxford University took up the problem again. They isolated penicillin in a form that was fairly pure (by standards then current) and demonstrated its potency and relative lack of toxicity. By then World War II had begun, and techniques to facilitate commercial production were developed in the United States. By 1944 adequate amounts were available to meet the extraordinary needs of wartime.

Antituberculous drugs. While penicillin is the most useful and the safest antibiotic, it suffers from certain disadvantages. The most important of these is that it is not active against *Mycobacterium tuberculosis*, the bacillus of tuberculosis. In view of the importance of tuberculosis as a public health hazard, this is a serious defect. The position was rapidly rectified when, in 1944, Selman A. Waksman and his colleagues announced the discovery of streptomycin from cultures of a soil organism, *Streptomyces griseus*, and stated that it was active against *M. tuberculosis*. Subsequent clinical trials amply confirmed this claim. Streptomycin suffers, however, from the great disadvantage that the tubercle bacillus tends to become resistant to it. Fortunately, other drugs became available

to supplement it, the two most important being para-aminosalicylic acid (or PAS) and isoniazid. With a combination of two or more of these preparations, the outlook in tuberculosis improved immeasurably. The disease was not conquered, but it was brought well under control.

Other antibiotics. Penicillin is not effective over the entire field of microorganisms pathogenic to humans. During the 1950s the search for antibiotics to fill this gap resulted in a steady stream of them, some with a much wider antibacterial range than penicillin (the so-called broad-spectrum antibiotics) and some capable of coping with those microorganisms that are inherently resistant to penicillin or that have developed resistance through exposure to penicillin.

This tendency of microorganisms to develop resistance to penicillin at one time threatened to become almost as serious a problem as the development of resistance to streptomycin by the bacillus of tuberculosis. Fortunately, early appreciation of the problem by clinicians resulted in more discriminate use of penicillin. Scientists continued to look for means of obtaining new varieties of penicillin, and their researches produced the so-called semisynthetic antibiotics, some of which are active when taken by mouth, while others are effective against microorganisms that have developed resistance to the earlier form of penicillin.

IMMUNOLOGY

Dramatic though they undoubtedly were, the advances in chemotherapy still left one important area vulnerable, that of the viruses. It was in bringing viruses under control that advances in immunology—the study of immunity—played such a striking part. One of the paradoxes of medicine is that the first large-scale immunization against a viral disease was instituted and established long before viruses were discovered. When Edward Jenner introduced vaccination against the virus that causes smallpox, the identification of viruses was still 100 years in the future. It took almost another half century to discover an effective method of producing antiviral vaccines that were both safe and effective.

In the meantime, however, the process by which the body reacts against infectious organisms to generate immunity became better understood. In Paris, Élie Metchnikoff had already detected the role of white blood cells in the immune reaction, and Jules Bordet had identified antibodies in the blood serum. The mechanisms of antibody activity were used to devise diagnostic tests for a number of diseases. In 1906 August von Wassermann gave his name to the blood test for syphilis, and in 1908 the tuberculin test—the skin test for tuberculosis—came into use. At the same time, methods of producing effective substances for inoculation were improved, and immunization against bacterial diseases made rapid progress.

Antibacterial vaccination. *Typhoid.* In 1897 the English bacteriologist Almroth Wright introduced a vaccine prepared from killed typhoid bacilli as a preventive of typhoid. Preliminary trials in the Indian army produced excellent results, and typhoid vaccination was adopted for the use of British troops serving in the South African War. Unfortunately, the method of administration was inadequately controlled, and the government sanctioned inoculations only for soldiers that “voluntarily presented themselves for this purpose prior to their embarkation for the seat of war.” The result was that, according to the official records, only 14,626 men volunteered out of a total strength of 328,244 who served during the three years of the war. Although later analysis showed that inoculation had had a beneficial effect, there were 57,684 cases of typhoid—approximately one in six of the British troops engaged—with 9,022 deaths.

A bitter controversy over the merits of the vaccine followed, but before the outbreak of World War I immunization had been officially adopted by the army. Comparative statistics would seem to provide striking confirmation of the value of antityphoid inoculation, even allowing for the better sanitary arrangements in the latter war. In the South African War the annual incidence of enteric infections (typhoid and paratyphoid) was 105 per 1,000 troops, and the annual death rate was 14.6 per 1,000; the comparable

Wasser-
mann
and the
tuberculin
test

Strepto-
mycin,
PAS, and
isoniazid

figures for World War I were 2.35 and 0.139, respectively.

It is perhaps a sign of the increasingly critical outlook that developed in medicine in the post-1945 era that experts continued to differ on some aspects of typhoid immunization. There was no question as to its fundamental efficacy, but there was considerable variation of opinion as to the best vaccine to use and the most effective way of administering it. Moreover, it was often difficult to decide to what extent the decline in typhoid was attributable to improved sanitary conditions and what to the greater use of the vaccine.

Tetanus. The other great hazard of war that was brought under control in World War I was tetanus. This was achieved by the prophylactic injection of tetanus antitoxin into all wounded men. The serum was originally prepared by the bacteriologists Emil von Behring and Shibasaburo Kitasato in 1890-92, and the results of this first large-scale trial amply confirmed its efficacy. (Tetanus antitoxin is a sterile solution of antibody globulins—a type of blood protein—from immunized horses or cattle.)

It was not until the 1930s, however, that an efficient vaccine, or toxoid, as it is known in the cases of tetanus and diphtheria, was produced against tetanus. (Tetanus toxoid is a preparation of the toxin—or poison—produced by the microorganism; injected into humans, it stimulates the body's own defenses against the disease, thus bringing about immunity.) Again, a war was to provide the opportunity for testing on a large scale, and experience with tetanus toxoid in World War II indicated that it gave a high degree of protection.

Diphtheria. The story of diphtheria is comparable to that of tetanus, though even more dramatic. First, as with tetanus antitoxin, came the preparation of diphtheria antitoxin by Behring and Kitasato in 1890. As the antitoxin came into general use for the treatment of cases, the death rate began to decline. There was no significant fall in the number of cases, however, until a toxin-antitoxin mixture, introduced by Behring in 1913, was used to immunize children. A more effective toxoid was introduced by the French bacteriologist Gaston Ramon in 1923, and with subsequent improvements this became one of the most effective vaccines available in medicine. Where mass immunization of children with the toxoid was practiced, as in the United States and Canada beginning in the late 1930s and in England and Wales in the early 1940s, cases of diphtheria and deaths from the disease became almost nonexistent. In England and Wales, for instance, the number of deaths fell from an annual average of 1,830 in 1940-44 to zero in 1969. Administration of a combined vaccine against diphtheria, pertussis (whooping cough), and tetanus (DPT) is recommended for young children. Although an increasing number of dangerous side effects from the DPT vaccine have been reported, it continues to be used in most countries because of the protection it affords.

BCG vaccine for tuberculosis. If, as is universally accepted, prevention is better than cure, immunization is the ideal way of dealing with diseases caused by microorganisms. An effective, safe vaccine protects the individual from disease, whereas chemotherapy merely copes with the infection once the individual has been affected. In spite of its undoubted value, however, immunization has been a recurring source of dispute. Like vaccination against typhoid (and against poliomyelitis later), tuberculosis immunization evoked widespread contention.

In 1908 Albert Calmette, a pupil of Pasteur, and Camille Guérin produced an avirulent (weakened) strain of the tubercle bacillus. About 13 years later, vaccination of children against tuberculosis was introduced, with a vaccine made from this avirulent strain and known as BCG (bacillus Calmette-Guérin) vaccine. Although it was adopted in France, Scandinavia, and elsewhere, British and U.S. authorities frowned upon its use on the grounds that it was not safe and that the statistical evidence in its favour was not convincing.

One of the stumbling blocks in the way of its widespread adoption was what came to be known as the Lübeck disaster. In the spring of 1930, 249 infants were vaccinated with BCG vaccine in Lübeck, Ger.; by autumn, 73 of

the 249 were dead. Criminal proceedings were instituted against those responsible for giving the vaccine. The final verdict was that the vaccine had been contaminated, and the BCG vaccine itself was exonerated from any responsibility for the deaths. A bitter controversy followed, but in the end the protagonists of the vaccine won when a further trial showed that the vaccine was safe and that it protected four out of five of those vaccinated.

Immunization against viral diseases. With the exception of smallpox, it was not until well into the 20th century that efficient viral vaccines became available. In fact, it was not until the 1930s that much began to be known about viruses. The two developments that contributed most to the rapid growth in knowledge after that time were the introduction of tissue culture as a means of growing viruses in the laboratory and the availability of the electron microscope. Once the virus could be cultivated with comparative ease in the laboratory, the research worker could study it with care and evolve methods for producing one of the two requirements for a safe and effective vaccine: either a virus that was so attenuated, or weakened, that it could not produce the disease for which it was responsible in its normally virulent form; or a killed virus that retained the faculty of inducing a protective antibody response in the vaccinated individual.

The first of the viral vaccines to result from these advances was for yellow fever, developed by the microbiologist Max Theiler in the late 1930s. About 1945 the first relatively effective vaccine was produced for influenza; in 1954 the American physician Jonas E. Salk introduced a vaccine for poliomyelitis; and in 1960 an oral poliomyelitis vaccine, developed by the virologist Albert B. Sabin, came into wide use.

These vaccines went far toward bringing under control three of the major diseases of the time although, in the case of influenza, a major complication is the disturbing proclivity of the virus to change its character from one epidemic to another. Even so, sufficient progress has been made to ensure that a pandemic like the one that swept the world in 1918-19, killing more than 15,000,000 people, is unlikely to occur again. Centres are now equipped to monitor outbreaks of influenza throughout the world in order to establish the identity of the responsible viruses and, if necessary, take steps to produce appropriate vaccines.

During the 1960s effective vaccines came into use for measles and rubella (German measles). Both evoked a certain amount of controversy. In the case of measles in the Western world it was contended that, if acquired in childhood, it is not a particularly hazardous malady, and the naturally acquired disease evokes permanent immunity in the vast majority of cases. Conversely, the vaccine induces a certain number of adverse reactions, and the duration of the immunity it produces is problematical. In the end the official view was that universal measles vaccination is to be commended.

The situation with rubella vaccination was different. This is a fundamentally mild affliction, and the only cause for anxiety is its proclivity to induce congenital deformities if a pregnant woman should acquire the disease. Once an effective vaccine was available, the problem was the extent to which it should be used. Ultimately the consensus was reached that all girls who had not already had the disease should be vaccinated at about 12 years. In the United States children are routinely immunized against measles, mumps, and rubella at the age of 15 months.

The immune response. With advances in cell biology in the second half of the 20th century came a more profound understanding of both normal and abnormal conditions in the body. Electron microscopy enabled observers to peer more deeply into the structures of the cell, and chemical investigations revealed clues to their functions in the cell's intricate metabolism. The overriding importance of the nuclear genetic material DNA (deoxyribonucleic acid) in regulating the cell's protein and enzyme production lines became evident. A clearer comprehension also emerged of the ways in which the cells of the body defend themselves by modifying their chemical activities to produce antibodies against injurious agents.

Up until the turn of the century, immunity referred

Poliomyelitis vaccine

Mass diphtheria immunization of children

Understanding the cell

mostly to the means of resistance of an animal to invasion by a parasite or microorganism. Around mid-century there arose a growing realization that immunity and immunology cover a much wider field and are concerned with mechanisms for preserving the integrity of the individual. The introduction of organ transplantation, with its dreaded complication of tissue rejection, brought this broader concept of immunology to the fore.

At the same time, research workers and clinicians began to appreciate the far-reaching implications of immunity in relation to endocrinology, genetics, tumour biology, and the biology of a number of other maladies. The so-called autoimmune diseases are caused by an aberrant series of immune responses by which the body's own cells are attacked. Suspicion is growing that a number of major disorders such as diabetes, rheumatoid arthritis, and multiple sclerosis may be caused by similar mechanisms.

In some conditions viruses invade the genetic material of cells and distort their metabolic processes. Such viruses may lie dormant for many years before becoming active. This may be the underlying cause of many cancers, in which cells escape from the usual constraints imposed upon them by the normal body. The dreaded affliction of acquired immune deficiency syndrome (AIDS) is caused by a virus that has a long dormant period and then attacks the cells that produce antibodies. The result is that the affected person is not able to generate an immune response to infections or malignancies.

AIDS

ENDOCRINOLOGY

At the beginning of the 20th century, endocrinology was in its infancy. Indeed, it was not until 1905 that Ernest H. Starling, one of the many brilliant pupils of Edward Sharpey-Schafer, the dean of British physiology during the early decades of the century, introduced the term hormone for the internal secretions of the endocrine glands. In 1891 the English physician George Redmayne Murray achieved the first success in treating myxedema (the common form of hypothyroidism) with an extract of the thyroid gland. Three years later, Sharpey-Schafer and George Oliver demonstrated in extracts of the adrenal glands a substance that raised the blood pressure; and in 1901 Jokichi Takamine, a Japanese chemist working in the United States, isolated this active principle, known as epinephrine or adrenaline.

Insulin. During the first two decades of the century, steady progress was made in the isolation, identification, and study of the active principles of the various endocrine glands, but the outstanding event of the early years was the discovery of insulin by Frederick Banting, Charles H. Best, and J.J.R. Macleod in 1921. Almost overnight the lot of the diabetic patient changed from a sentence of almost certain death to a prospect not only of survival but of a long and healthy life.

For more than 30 years, some of the greatest minds in physiology had been seeking the cause of diabetes mellitus. In 1889 the German physicians Joseph von Mering and Oskar Minkowski had shown that removal of the pancreas in dogs produced the disease. In 1901 the American pathologist Eugene L. Opie described degenerative changes in the clumps of cells in the pancreas known as the islets of Langerhans, thus confirming the association between failure in the function of these cells and diabetes. Sharpey-Schafer concluded that the islets of Langerhans secrete a substance that controls the metabolism of carbohydrate. Then Banting, Best, and Macleod, working at the University of Toronto, succeeded in isolating the elusive hormone and gave it the name insulin.

Insulin was available in a variety of forms, but synthesis on a commercial scale was not achieved, and the only source of the hormone was the pancreas of animals. One of its practical disadvantages is that it has to be given by injection; consequently an intense search was conducted for some alternative substance that would be active when taken by mouth. Various preparations—oral hypoglycemic agents, as they are known—appeared that were effective to a certain extent in controlling diabetes, but evidence indicated that these were only of value in relatively mild cases of the disease. For the person with advanced dia-

betes, a normal, healthy life remained dependent upon the continuing use of insulin injections.

Cortisone. Another major advance in endocrinology came from the Mayo Clinic, in Rochester, Minn. In 1949 Philip S. Hench and his colleagues announced that a substance isolated from the cortex of the adrenal gland had a dramatic effect upon rheumatoid arthritis. This was compound E, or cortisone, as it came to be known, which had been isolated by Edward C. Kendall in 1935. Cortisone and its many derivatives proved to be potent as anti-inflammatory agents. Although it is not a cure for rheumatoid arthritis, as a temporary measure cortisone can often control the acute exacerbation caused by the disease and can provide relief in other conditions, such as acute rheumatic fever, certain kidney diseases, certain serious diseases of the skin, and some allergic conditions, including acute exacerbations of asthma. Of even more long-term importance is the valuable role it has as a research tool.

Sex hormones. Not the least of the advances in endocrinology was the increasing knowledge and understanding of the sex hormones. This culminated in the application of this knowledge to the problem of birth control. After an initial stage of hesitancy, the contraceptive pill, with its basic rationale of preventing ovulation, was accepted by the vast majority of family-planning organizations and many gynecologists as the most satisfactory method of contraception. Its risks, practical and theoretical, introduced a note of caution, but this was not sufficient to detract from the wide appeal induced by its effectiveness and ease of use.

The contraceptive pill

VITAMINS

In the field of nutrition, the outstanding advance of the 20th century was the discovery and the appreciation of the importance to health of the "accessory food factors," or vitamins. Various workers had shown that animals did not thrive on a synthetic diet containing all the correct amounts of protein, fat, and carbohydrate; they even suggested that there must be some unknown ingredients in natural food that were essential for growth and the maintenance of health. But little progress was made in this field until the classical experiments of the English biologist F. Gowland Hopkins were published in 1912. These were so conclusive that there could be no doubt that what he termed "accessory substances" were essential for health and growth.

The name vitamin was suggested for these substances by the biochemist Casimir Funk in the belief that they were amines, certain compounds derived from ammonia. In due course, when it was realized that they were not amines, the term was altered to vitamin.

Once the concept of vitamins was established on a firm scientific basis it was not long before their identity began to be revealed. Soon there was a long series of vitamins, best known by the letters of the alphabet after which they were originally named when their chemical identity was still unknown. By supplementing the diet with foods containing particular vitamins, deficiency diseases such as rickets (due to deficiency of vitamin D) and scurvy (due to lack of vitamin C, or ascorbic acid) practically disappeared from Western countries, while deficiency diseases such as beriberi (caused by lack of vitamin B₁, or thiamine), which were endemic in Eastern countries, either disappeared or could be remedied with the greatest of ease.

The isolation of vitamin B₁₂, or cyanocobalamin, was of particular interest because it almost rounded off the fascinating story of how pernicious anemia was brought under control. Throughout the first two decades of the century, the diagnosis of pernicious anemia, like that of diabetes mellitus, was nearly equivalent to a death sentence. Unlike the more common form of so-called secondary anemia, it did not respond to the administration of suitable iron salts, and no other form of treatment touched it; hence, the grimly appropriate title of pernicious anemia.

In the early 1920s, George R. Minot, one of the many brilliant investigators that Harvard University has contributed to medical research, became interested in work being done by the American pathologist George H. Whip-

ple on the beneficial effects of raw beef liver in severe experimental anemia. With a Harvard colleague, William P. Murphy, he decided to investigate the effect of raw liver in patients with pernicious anemia, and in 1926 they were able to announce that this form of therapy was successful. The validity of their findings was amply confirmed, and the fear of pernicious anemia came to an end.

Liver
therapy in
pernicious
anemia

As so often happens in medicine, many years were to pass before the rationale of liver therapy in pernicious anemia was fully understood. In 1948, however, almost simultaneously in the United States and Britain, the active principle, cyanocobalamin, was isolated from liver, and this vitamin became the standard treatment for pernicious anemia.

MALIGNANT DISEASE

While progress was the hallmark of medicine after the beginning of the 20th century, there is one field in which a gloomier picture must be painted, that of malignant disease, or cancer. It is the second most common cause of death in most Western countries in the second half of the 20th century, being exceeded only by deaths from heart disease. Some progress, however, has been achieved. The causes of the various types of malignancies are not known, but many more methods are available for attacking the problem; surgery remains the principal therapeutic standby, but radiotherapy and chemotherapy are increasingly used.

Soon after the discovery of radium was announced, in 1898, its potentialities in treating cancer were realized; in due course it assumed an important role in therapy. Simultaneously, deep X-ray therapy was developed, and with the atomic age came the use of radioactive isotopes. (A radioactive isotope is an unstable variant of a substance that has a stable form; during the process of breaking down, the unstable form emits radiation.) High-voltage X-ray therapy and radioactive isotopes have largely replaced radium. Whereas irradiation long depended upon X rays generated at 250 kilovolts, machines that are capable of producing X rays generated at 8,000 kilovolts and betatrons of up to 22,000,000 electron volts (MeV) have come into clinical use.

The most effective of the isotopes is radioactive cobalt. Telecobalt machines (those that hold the cobalt at a distance from the body) are available containing 2,000 curies or more of the isotope, an amount equivalent to 3,000 grams of radium and sending out a beam equivalent to that from a 3,000-kilovolt X-ray machine.

Of even more significance have been the developments in the chemotherapy of cancer. Nothing remotely resembling a chemotherapeutic cure has been achieved, but in certain forms of malignant disease, such as leukemia, which cannot be treated by surgery, palliative effects have been achieved that prolong life and allow the patient in many instances to lead a comparatively normal existence.

Chemo-
therapy
for cancer

Fundamentally, however, perhaps the most important advance of all in this field has been the increasing appreciation of the importance of prevention. The discovery of the relationship between cigarette smoking and lung cancer is the classic example. Less publicized, but of equal import, is the continuing supervision of new techniques in industry and food manufacture in an attempt to ensure that they do not involve the use of cancer-causing substances.

TROPICAL MEDICINE

The first half of the 20th century witnessed the virtual conquest of three of the major diseases of the tropics: malaria, yellow fever, and leprosy. At the turn of the century, as for the preceding two centuries, quinine was the only known drug to have any appreciable effect on malaria. With the increasing development of tropical countries and rising standards of public health, it became obvious that quinine was not completely satisfactory. Intensive research between World Wars I and II indicated that several synthetic compounds were more effective. The first of these to become available, in 1934, was quinacrine (known as mepacrine, Atabrine, or Atebrin). In World War II it amply fulfilled the highest expectations and helped to reduce disease among Allied troops in Africa, Southeast Asia,

and the Far East. A number of other effective antimalarial drugs subsequently became available.

An even brighter prospect—the virtual eradication of malaria—was opened up by the introduction, during World War II, of the insecticide DDT (1,1,1-trichloro-2,2,-bis[*p*-chlorophenyl]ethane, or dichlorodiphenyltrichloroethane). It had long been realized that the only effective way of controlling malaria was to eradicate the anopheline mosquitoes that transmit the disease. Older methods of mosquito control, however, were cumbersome and expensive. The lethal effect of DDT on the mosquito, its relative cheapness, and its ease of use on a widespread scale provided the answer. An intensive worldwide campaign, sponsored by the World Health Organization, was planned and went far toward bringing malaria under control.

The major problem encountered with respect to effectiveness was that the mosquitoes were able to develop a resistance to DDT; but the introduction of other insecticides, such as dieldrin and lindane (BHC), helped to overcome this difficulty. In recent years the use of these and other insecticides has been strongly criticized by ecologists, however.

Yellow fever is another mosquito-transmitted disease, and the prophylactic value of modern insecticides in its control was almost as great as in the case of malaria. The forest reservoirs of the virus present a more difficult problem, but the combined use of immunization and insecticides did much to bring this disease under control.

Until the 1940s the only drugs available for treating leprosy were the chaulmoogra oils and their derivatives. These, though helpful, were far from satisfactory. In the 1940s the group of drugs known as the sulfones appeared, and it soon became apparent that they were infinitely better than any other group of drugs in the treatment of leprosy. Several other drugs later proved promising. Although there is as yet no known cure—in the strict sense of the term—for leprosy, the outlook has so changed that there are good grounds for believing that this age-old scourge can be brought under control and the victims of the disease saved from those dreaded mutilations that have given leprosy such a fearsome reputation throughout the ages. (Wm.A.R.T./P.Rh.)

Control of
leprosy

Surgery in the 20th century

THE OPENING PHASE

Three seemingly insuperable obstacles beset the surgeon in the years before the mid-19th century: pain, infection, and shock. Once these were overcome, the surgeon believed that he could burst the bonds of centuries and become the master of his craft. There is more, however, to anesthesia than putting the patient to sleep. Infection, despite first antiseptics (destruction of microorganisms present) and later asepsis (avoidance of contamination), is still an ever-present menace; and shock continues to perplex physicians. But in the 20th century, surgery has progressed farther, faster, and more dramatically than in all preceding ages.

The situation encountered. The shape of surgery that entered the new century was clearly recognizable as the forerunner of today's, blurred and hazy though the outlines may now seem. The operating theatre still retained an aura of the past, when the surgeon played to his audience and the patient was little more than a stage prop. In most hospitals it was a high room lit by a skylight, with tiers of benches rising above the narrow, wooden operating table. The instruments, kept in glazed or wooden cupboards around the walls, were of forged steel, unplated, and with handles of wood or ivory.

The means to combat infection hovered between antiseptics and asepsis. Instruments and dressings were mostly sterilized by soaking them in dilute carbolic acid (or other antiseptic), and the surgeon often endured a gown freshly wrung out in the same solution. Asepsis gained ground fast, however. It had been born in the Berlin clinic of Ernst von Bergmann where, in 1886, steam sterilization had been introduced. Gradually, this led to the complete aseptic ritual, which has as its basis the bacterial cleanliness (as opposed to social cleanliness) of everything that

comes in contact with the wound. Hermann Kümmell, of Hamburg, devised the routine of "scrubbing up." In 1890 William Stewart Halsted, of Johns Hopkins University, had rubber gloves specially made for operating, and in 1896 Johannes von Mikulicz-Radecki, a Pole working at Breslau, Ger., invented the gauze mask.

Many surgeons, brought up in a confused misunderstanding of the antiseptic principle—believing that carbolic would cover a multitude of sins, many of which they were ignorant of committing—failed to grasp what asepsis was all about. Thomas Annandale, for example, blew through his catheters to make sure that they were clear, and many an instrument, dropped accidentally, was simply given a quick wipe and returned to use. Tradition died hard, and asepsis had an uphill struggle before it was fully accepted. "I believe firmly that more patients have died from the use of gloves than have ever been saved from infection by their use," wrote W.P. Carr, an American, in 1911. Over the years, however, a sound technique was evolved as the foundation for the growth of modern surgery.

Anesthesia at the turn of the century Anesthesia, at the turn of the century, progressed slowly. Few physicians made a career of the subject, and frequently the patient was rendered unconscious by a student, a nurse, or a porter wielding a rag and bottle. Chloroform was overwhelmingly more popular than ether, on account of its ease of administration, despite the fact that it was liable to kill by stopping the heart.

Although by the end of the first decade, nitrous oxide (laughing gas) combined with ether had displaced—but by no means entirely—the use of chloroform, the surgical problems were far from ended. For years to come the abdominal surgeon besought the anesthetist to deepen the level of anesthesia and thus relax the abdominal muscles; the anesthetist responded to the best of his ability, acutely aware that the deeper he went, the closer the patient was to death. When other anesthetic agents were discovered, the anesthetist came into his own, and many advances in spheres such as brain and heart surgery would have been impossible without his skill.

The third obstacle, shock, is perhaps the most complex and the most difficult to define satisfactorily. The only major cause properly appreciated at the start of the 20th century was loss of blood, and once that had occurred nothing, in those days, could be done. And so, the study of shock—its causes, its effects on human physiology, and its prevention and treatment—became all-important to the progress of surgery.

In the latter part of the 19th century, then, surgeons had been liberated from the age-old bogies of pain, pus, and hospital gangrene. Hitherto, operations had been restricted to amputations, cutting for stone in the bladder, tying off arterial aneurysms (bulging and thinning of artery walls), repairing hernias, and a variety of procedures that could be done without going too deeply beneath the skin. But the anatomical knowledge, a crude skill derived from practice on dead bodies, and above all the enthusiasm, were there waiting. Largely ignoring the mass of problems they uncovered, surgeons launched forth into an exploration of the human body.

They acquired a reputation for showmanship; but much of their surgery, though speedy and spectacular, was rough and ready. There were a few who developed supreme skill and dexterity and could have undertaken a modern operation with but little practice; indeed, some devised the very operations still in use today. One such was Theodor Billroth, head of the surgical clinic at Vienna, who collected a formidable list of successful "first" operations. He represented the best of his generation—a surgical genius, an accomplished musician, and a kind, gentle man who brought the breath of humanity to his work. Moreover, the men he trained, including von Mikulicz, Vincenz Czerny, and Anton von Eiselsberg, consolidated the brilliant start that he had given to abdominal surgery in Europe.

Billroth's influence

Changes before World War I. The opening decade of the 20th century was a period of transition. Flamboyant exhibitionism was falling from favour as surgeons, through experience, learned the merits of painstaking, conscientious operation—treating the tissues gently and carefully controlling every bleeding point. The individualist was not

submerged, however, and for many years the development of the various branches of surgery rested on the shoulders of a few clearly identifiable men. Teamwork on a large scale arrived only after World War II. The surgeon, at first, was undisputed master in his own wards and theatre. But as time went on and he found he could not solve his problems alone, he called for help from specialists in other fields of medicine and, even more significantly, from his colleagues in other scientific disciplines.

The increasing scope of surgery led to specialization. Admittedly, most general surgeons had a special interest, and for a long time there had been an element of specialization in such fields as ophthalmology, orthopedics, obstetrics, and gynecology; but before long it became apparent that, to achieve progress in certain areas, surgeons had to concentrate their attention on that particular subject.

Abdominal surgery. By the start of the 20th century, abdominal surgery, which provided the general surgeon with the bulk of his work, had grown beyond infancy, thanks largely to Billroth. In 1881 he had performed the first successful removal of part of the stomach for cancer. His next two cases were failures, and he was stoned in the streets of Vienna. Yet, he persisted and by 1891 had carried out 41 more of these operations with 16 deaths—a remarkable achievement for that era.

Peptic ulcers (gastric and duodenal) appeared on the surgical scene (perhaps as a new disease, but more probably because they had not been diagnosed previously), and in 1881 Ludwig Rydygier cured a young woman of her gastric ulcer by removing it. Bypass operations—gastroenterostomies—soon became more popular, however, and enjoyed a vogue that lasted into the 1930s, even though fresh ulcers at the site of the juncture were not uncommon.

The other end of the alimentary tract was also subjected to surgical intervention; cancers were removed from the large bowel and rectum with mortality rates that gradually fell from 80 to 60 to 20 to 12 percent as the surgeons developed their skill. In 1908 the British surgeon Ernest Miles carried out the first abdominoperineal resection for cancer of the rectum; that is, the cancer was attacked both from the abdomen and from below through the perineum (the area between the anus and the genitals), either by one surgeon, who actually did two operations, or by two working together. This technique formed the basis for all future developments.

Much of the new surgery in the abdomen was for cancer, but not all. Appendectomy became the accepted treatment for appendicitis (in appropriate cases) in the United States before the close of the 19th century; but in Great Britain surgeons were reluctant to remove the organ until 1902, when King Edward VII's coronation was dramatically postponed on account of his appendicitis. The publicity attached to his operation caused the disease and its surgical treatment to become fashionable—despite the fact that the royal appendix remained in the King's abdomen; the surgeon, Frederic Treves, had merely drained the abscess.

Appendectomy

Neurosurgery. Though probably the most demanding of all the surgical specialties, neurosurgery was nevertheless one of the first to emerge. The techniques and principles of general surgery were inadequate for work in such a delicate field. William Macewen, a Scottish general surgeon of outstanding versatility, and Victor Alexander Haden Horsley, the first British neurosurgeon, showed that the surgeon had much to offer in the treatment of disease of the brain and spinal cord. Macewen, in 1893, recorded 19 patients operated on for brain abscess, 18 of whom were cured; at that time most other surgeons had 100 percent mortality rates for the condition. His achievement remained unequaled until the discovery of penicillin.

An American, Harvey Williams Cushing, almost by himself consolidated neurosurgery as a specialty. From 1905 on, he advanced neurosurgery through a series of operations and through his writings. Tumours, epilepsy, trigeminal neuralgia, and pituitary disorders were among the conditions he treated successfully.

Radiology. In 1895 a development at the University of Würzburg had far-reaching effects on medicine and surgery, opening up an entirely fresh field of the diagnosis and study of disease and leading to a new form of

treatment, radiation therapy. This was the discovery of X rays by Wilhelm Conrad Röntgen, a professor of physics. Within months of the discovery there was an extensive literature on the subject: Robert Jones, a British surgeon, had localized a bullet in a boy's wrist before operating; stones in the urinary bladder and gallbladder had been demonstrated; and fractures had been displayed.

X rays in research

Experiments began on introducing substances that are opaque to X rays into the body to reveal organs and formations, both normal and abnormal. Walter Cannon, a Boston physiologist, used X rays in 1898 in his studies of the alimentary tract. Friedrich Voelcker, of Heidelberg, devised retrograde pyelography (introduction of the radiopaque medium into the kidney pelvis by way of the ureter) for the study of the urinary tract in 1905; in Paris in 1921, Jean Sicard X-rayed the spinal canal with the help of an oily iodine substance, and the next year he did the same for the bronchial tree; and in 1924 Everts Graham, of St. Louis, used a radiopaque contrast medium to view the gallbladder. Air was also used to provide contrast; in 1918, at Johns Hopkins, Walter Dandy injected air into the ventricles (liquid-filled cavities) of the brain.

The problems of injecting contrast media into the blood vessels took longer to solve, and it was not until 1927 that António Moniz, of Lisbon, succeeded in obtaining pictures of the arteries of the brain. Eleven years later, George Robb and Israel Steinberg of New York overcame some of the difficulties of cardiac catheterization (introduction of a small tube into the heart by way of veins or arteries) and were able to visualize the chambers of the heart on X-ray film. After much research, a further refinement came in 1962, when Frank Sones and Earl K. Shirey of Cleveland showed how to introduce the contrast medium into the coronary arteries.

WORLD WAR I

The battlefields of the 20th century stimulated the progress of surgery and taught the surgeon innumerable lessons, which were subsequently applied in civilian practice. Regrettably, though, the principles of military surgery and casualty evacuation, which can be traced back to the Napoleonic wars, had to be learned over again.

Treatment of contaminated wounds

World War I broke, quite dramatically, the existing surgical hierarchy and rule of tradition. No longer did the European surgeon have to waste his best years in apprenticeship before seating himself in his master's chair. Suddenly, young surgeons in the armed forces began confronting problems that would have daunted their elders. Furthermore, their training had been in "clean" surgery performed under aseptic conditions. Now they found themselves faced with the need to treat large numbers of grossly contaminated wounds in improvised theatres. They rediscovered debridement (the surgical excision of dead and dying tissue and the removal of foreign matter).

The older surgeons cried "back to Lister," but antiseptics, no matter how strong, were no match for putrefaction and gangrene. One method of antiseptic irrigation—devised by Alexis Carrel and Henry Dakin and called the Carrel-Dakin treatment—was, however, beneficial, but only after the wound had been adequately debrided. The scourges of tetanus and gas gangrene were controlled to a large extent by antitoxin and antiserum injections, yet surgical treatment of the wound remained an essential requirement.

Abdominal casualties fared badly for the first year of the war, because experience in the utterly different circumstances of the South African War had led to a belief that these men were better left alone surgically. Fortunately, the error of continuing with such a policy 15 years later was soon appreciated, and every effort was made to deliver the wounded men to a suitable surgical unit with all speed. Little progress was made with chest wounds beyond opening up the wound even further to drain pus from the pleural cavity between the chest wall and the lungs.

Perhaps the most worthwhile and enduring benefit to flow from World War I was rehabilitation. For almost the first time, surgeons realized that their work did not end with a healed wound. In 1915 Robert Jones set up special facilities for orthopedic patients, and at about the same time Harold Gillies founded British plastic surgery

in a hut at Sidcup, Kent. In 1917 Gillies popularized the pedicle type of skin graft (the type of graft in which skin and subcutaneous tissue are left temporarily attached for nourishment to the site from which the graft was taken). Since then plastic surgery has given many techniques and principles to other branches of surgery.

BETWEEN THE WORLD WARS

The years between the two world wars may conveniently be regarded as the time when surgery consolidated its position. A surprising number of surgical firsts and an amazing amount of fundamental research had been achieved even in the late 19th century, but the knowledge and experience could not be converted to practical use because the human body could not survive the onslaught. In the years between World Wars I and II, it was realized that physiology—in its widest sense, including biochemistry and fluid and electrolyte balance—was of major importance along with anatomy, pathology, and surgical technique.

The problem of shock. The first problem to be tackled was shock, which was, in brief, found to be due to a decrease in the effective volume of the circulation. To combat shock, the volume had to be restored, and the obvious substance was blood itself. In 1901 Karl Landsteiner, then in Austria, discovered the ABO blood groups, and in 1914 sodium citrate was added to freshly drawn blood to prevent clotting. Blood was occasionally transfused during World War I, but three-quarters of a pint was considered a large amount. These transfusions were given by directly linking the vein of a donor with that of the recipient. The continuous drip method, in which blood flows from a flask, was introduced by Hugh Marriott and Alan Kekwick at the Middlesex Hospital, London, in 1935.

Blood transfusion

As blood transfusions increased in frequency and volume, blood banks were required. Although it took another world war before these were organized on a large scale, the first tentative steps were taken by Sergey Sergeevich Yudin, of Moscow, who, in 1933, used cadaver blood, and by Bernard Fantus, of Chicago, who, four years later, used living donors as his source of supply. Saline solution, plasma, artificial plasma expanders, and other solutions are now also used in the appropriate circumstances.

Sometimes after operations (especially abdominal operations), the gut becomes paralyzed. It is distended, and quantities of fluid pour into it, dehydrating the body. In 1932 Owen Wangenstein, at the University of Minnesota, advised decompressing the bowel, and in 1934 two other Americans, Thomas Miller and William Abbott, of Philadelphia, invented an apparatus for this purpose, a tube with an inflatable balloon on the end that could be passed into the small intestine. The fluid lost from the tissues was replaced by a continuous intravenous drip of saline solution on the principle described by Rudolph Matas, of New Orleans, in 1924. These techniques dramatically improved abdominal surgery, especially in cases of obstruction, peritonitis (inflammation of the abdominal membranes), and acute emergencies generally, since they made it possible to keep the bowel empty and at rest.

Anesthesia and thoracic surgery. The strides taken in anesthesia from the 1920s onward allowed surgeons much more freedom. Rectal anesthesia had never proved satisfactory, and the first improvement on the combination of nitrous oxide, oxygen, and ether was the introduction of the general anesthetic cyclopropane by Ralph Waters of Madison, Wis., in 1933. Soon afterward, intravenous anesthesia was introduced; John Lundy of the Mayo Clinic brought to a climax a long series of trials by many workers when he used Pentothal (thiopental sodium, a barbiturate) to put a patient peacefully to sleep. Then, in 1942, Harold Griffith and G. Enid Johnson, of Montreal, produced muscular paralysis by the injection of a purified preparation of curare. This was harmless since, by then, the anesthetist was able to control the patient's respiration.

If there was one person who was aided more than any other by the progress in anesthesia, it was the thoracic (chest) surgeon. What had bothered him previously was the collapse of the lung, which occurred whenever the pleural cavity was opened. Since the end of the 19th century, many and ingenious methods had been devised

Preventing lung collapse to prevent this from happening. The best known was the negative pressure cabinet of Ernst Ferdinand Sauerbruch, then at Mikulicz' clinic at Breslau; the cabinet was first demonstrated in 1904 but was destined soon to become obsolete.

The solution lay in inhalational anesthesia administered under pressure. Indeed, when Théodore Tuffier, in 1891, successfully removed the apex of a lung for tuberculosis, this was the technique that he used; he even added an inflatable cuff around the tube inserted in the trachea to ensure a gas-tight fit. Tuffier was ahead of his time, however, and other surgeons and research workers wandered into confused and complex byways before Ivan Magill and Edgar Rowbotham, working at Gillies' plastic-surgery unit, found their way back to the simplicity of the endotracheal tube and positive pressure. In 1931 Ralph Waters showed that respiration could be controlled either by squeezing the anesthetic bag by hand or by using a small motor.

These advances allowed thoracic surgery to move into modern times. In the 1920s, operations had been performed mostly for infective conditions and as a last resort. The operations necessarily were unambitious and confined to collapse therapy, including thoracoplasty (removal of ribs), apicolysis (collapse of a lung apex and artificially filling the space), and phrenic crush (which paralyzed the diaphragm on the chosen side); to isolation of the area of lung to be removed by first creating pleural adhesions; and to drainage.

The technical problems of surgery within the chest were daunting until Harold Brunn of San Francisco reported six lobectomies (removals of lung lobes) for bronchiectasis with only one death. (In bronchiectasis one or more bronchi or bronchioles are chronically dilated and inflamed, with copious discharge of mucus mixed with pus.) The secret of Brunn's success was the use of intermittent suction after surgery to keep the cavity free of secretions until the remaining lobes of the lung could expand to fill the space. In 1931 Rudolf Nissen, in Berlin, removed an entire lung from a girl with bronchiectasis. She recovered to prove that the risks were not as bad as had been feared.

Cancer of the lung has become a major disease of the 20th century; perhaps it has genuinely increased, or perhaps modern techniques of diagnosis reveal it more often. As far back as 1913 a Welshman, Hugh Davies, removed a lower lobe for cancer, but a new era began when Everts Graham removed a whole lung for cancer in 1933. The patient, a doctor, was still alive at the time of Graham's death in 1957.

The thoracic part of the esophagus is particularly difficult to reach, but in 1909 the British surgeon Arthur Evans successfully operated on it for cancer. But results were generally poor until, in 1944, John Garlock, of New York, showed that it is possible to excise the esophagus and to bring the stomach up through the chest and join it to the pharynx. Lengths of colon are also used as grafts to bridge the gap.

WORLD WAR II AND AFTER

Once the principles of military surgery were relearned and applied to modern warfare, instances of death, deformity, and loss of limb were reduced to levels previously unattainable. This was due largely to a thorough reorganization of the surgical services, adapting them to prevailing conditions, so that casualties received the appropriate treatment at the earliest possible moment. Evacuation by air (first used in World War I) helped greatly in this respect. Diagnostic facilities were improved, and progress in anesthesia kept pace with the surgeon's demands. Blood was transfused in adequate—and hitherto unthinkable—quantities, and the blood transfusion service as it is known today came into being.

Surgical specialization and teamwork reached new heights with the creation of units to deal with the special problems of injuries to different parts of the body. But the most revolutionary change was in the approach to wound infections brought about by the use of sulfonamides and (after 1941) of penicillin. The fact that these drugs could never replace meticulous wound surgery was, however, another lesson learned only in the bitter school of experience.

When the war ended, surgeons returned to civilian life feeling that they were at the start of a completely new, exciting era; and indeed they were, for the intense stimulation of the war years had led to developments in many branches of science that could now be applied to surgery. Nevertheless, it must be remembered that these developments merely allowed surgeons to realize the dreams of their fathers and grandfathers; they opened up remarkably few original avenues. The two outstanding phenomena of the 1950s and 1960s—heart surgery and organ transplantation—both originated in a real and practical manner at the turn of the century.

Support from other technologies. At first, perhaps, the surgeon tried to do too much himself, but before long his failures taught him to share his problems with experts in other fields. This was especially so with respect to difficulties of biomedical engineering and the exploitation of new materials. The relative protection from infection given by antibiotics and chemotherapy allowed the surgeon to become far more adventurous than hitherto in repairing and replacing damaged or worn-out tissues with foreign materials. Much research was still needed to find the best material for a particular purpose and to make sure that it would be acceptable to the body.

Plastics, in their seemingly infinite variety, have come to be used for almost everything from suture material to heart valves; for strengthening the repair of hernias; for replacement of the head of the femur (first done by the French surgeon Jean Judet and his brother Robert-Louis Judet in 1950); for replacement of the lens of the eye after extraction of the natural lens for cataract; for valves to drain fluid from the brain in patients with hydrocephalus; and for many other applications. This is a far cry, indeed, from the unsatisfactory use of celluloid to restore bony defects of the face by the German surgeon Fritz Berndt in the 1890s. Inert metals, such as vitallium, have also found a place in surgery, largely in orthopedics for the repair of fractures and the replacement of joints.

The scope of surgery was further expanded by the introduction of the operating microscope. This brought the benefit of magnification particularly to neurosurgery and to ear surgery. In the latter it opened up a whole field of operations on the eardrum and within the middle ear. The principles of these operations were stated in 1951 and 1952 by two German surgeons, Fritz Zöllner and Horst Wullstein; and in 1952 Samuel Rosen of New York mobilized the footplate of the stapes to restore hearing in otosclerosis—a procedure attempted by the German Jean Kessel in 1876.

Although surgeons aim to preserve as much of the body as disease permits, they are sometimes forced to take radical measures to save life; when, for instance, cancer affects the pelvic organs. Pelvic exenteration (surgical removal of the pelvic organs and nearby structures) in two stages was devised by Allen Whipple of New York City, in 1935, and in one stage by Alexander Brunschwig, of Chicago, in 1937. Then, in 1960, Charles S. Kennedy, of Detroit, after a long discussion with Brunschwig, put into practice an operation that he had been considering for 12 years: hemicorporectomy—surgical removal of the lower part of the body. The patient died on the 11th day. The first successful hemicorporectomy (at the level between the lowest lumbar vertebra and the sacrum) was performed 18 months later by J. Bradley Aust and Karel B. Absolon, of Minnesota. This operation would never have been possible without all the technical, supportive, and rehabilitative resources of modern medicine.

Heart surgery. The attitude of the medical profession toward heart surgery was for long overshadowed by doubt and disbelief. Wounds of the heart could be sutured (first done successfully by Ludwig Rehn, of Frankfurt am Main, in 1896); the pericardial cavity—the cavity formed by the sac enclosing the heart—could be drained in purulent infections (as had been done by Larrey in 1824); and the pericardium could be partially excised for constrictive pericarditis when it was inflamed and constricted the movement of the heart (this operation was performed by Rehn and Sauerbruch in 1913). But little beyond these procedures found acceptance.

Use of plastics in surgery

Heart operations of the early 20th century

Yet, in the first two decades of the 20th century, much experimental work had been carried out, notably by the French surgeons Théodore Tuffier and Alexis Carrel. Tuffier, in 1912, operated successfully on the aortic valve. In 1923 Elliott Cutler of Boston used a tenotome, a tendon-cutting instrument, to relieve a girl's mitral stenosis (a narrowing of the mitral valve between the upper and lower chambers of the left side of the heart) and in 1925, in London, Henry Souttar used a finger to dilate a mitral valve in a manner that was 25 years ahead of its time. Despite these achievements, there was too much experimental failure, and heart disease remained a medical, rather than surgical, matter.

Resistance began to crumble in 1938, when Robert Gross successfully tied off a persistent ductus arteriosus (a fetal blood vessel between the pulmonary artery and the aorta). It was finally swept aside in World War II by the remarkable record of Dwight Harken, who removed 134 missiles from the chest—13 in the heart chambers—without the loss of one patient.

After the war, advances came rapidly, with the initial emphasis on the correction or amelioration of congenital defects. Gordon Murray, of Toronto, made full use of his amazing technical ingenuity to devise and perform many pioneering operations. And Charles Bailey of Philadelphia, adopting a more orthodox approach, was responsible for establishing numerous basic principles in the growing specialty.

Until 1953, however, the techniques all had one great disadvantage: they were done "blind." The surgeon's dream was to stop the heart so that he could see what he was doing and be allowed more time in which to do it. In 1952 this dream began to come true when Floyd Lewis, of Minnesota, reduced the temperature of the body so as to lessen its need for oxygen while he closed a hole between the two upper heart chambers, the atria. The next year John Gibbon, Jr., of Philadelphia brought to fulfillment the research he had begun in 1937; he used his heart-lung machine to supply oxygen while he closed a hole in the septum between the atria.

Unfortunately, neither method alone was ideal, but intensive research and development led, in the early 1960s, to their being combined as extracorporeal cooling. That is, the blood circulated through a machine outside the body, which cooled it (and, after the operation, warmed it); the cooled blood lowered the temperature of the whole body. With the heart dry and motionless, the surgeon operated on the coronary arteries; he inserted plastic patches over holes; he sometimes almost remodeled the inside of the heart. But when it came to replacing valves destroyed

by disease, he was faced with a difficult choice between human tissue and man-made valves, or even valves from animal sources.

Organ transplantation. In 1967 surgery arrived at a climax that made the whole world aware of its medico-surgical responsibilities when the South African surgeon Christiaan Barnard transplanted the first human heart. Reaction, both medical and lay, contained more than an element of hysteria. Yet, in 1964, James Hardy, of the University of Mississippi, had transplanted a chimpanzee's heart into a man; and in that year two prominent research workers, Richard Lower and Norman E. Shumway, had written: "Perhaps the cardiac surgeon should pause while society becomes accustomed to resurrection of the mythological chimera." Research had been remorselessly leading up to just such an operation ever since Charles Guthrie and Alexis Carrel, at the University of Chicago, perfected the suturing of blood vessels in 1905 and then carried out experiments in the transplantation of many organs, including the heart.

New developments in immunosuppression (the use of drugs to prevent organ rejection) have advanced the field of transplantation enormously. Kidney transplantation is now a routine procedure that is supplemented by dialysis with an artificial kidney (invented by Willem Kolff in wartime Holland) before and after the operation; mortality has been reduced to about 10 percent per year. Rejection of the transplanted heart by the patient's immune system was overcome to some degree in the 1980s with the introduction of the immunosuppressant cyclosporine; records show that many patients have lived for five or more years after the transplant operation.

The complexity of the liver and the unavailability of supplemental therapies such as the artificial kidney have contributed to the slow progress in liver transplantation (first performed in 1963 by Thomas Starzl). An increasing number of patients, especially children, have undergone successful transplantation; however, a substantial number may require retransplantation due to the failure of the first graft.

Lung transplants (first performed by Hardy in 1963) are difficult procedures, and much progress is yet to be made in preventing rejection. A combined heart-lung transplant is still in the experimental stage, but it is being met with increasing success; two-thirds of those receiving transplants are surviving, although complications such as infection are still common. Transplantation of all or part of the pancreas is not completely successful, and further refinements of the procedures (first performed in 1966 by Richard Lillehei) are needed. (R.G.R./Ed.)

THE PRACTICE OF MODERN MEDICINE

Health care and its delivery

The World Health Organization at its 1978 international conference held in the Soviet Union produced the Alma-Ata Health Declaration, which was designed to serve governments as a basis for planning health care that would reach people at all levels of society. The declaration reaffirmed that "health, which is a state of complete physical, mental and social well-being, and not merely the absence of disease or infirmity, is a fundamental human right and that the attainment of the highest possible level of health is a most important world-wide social goal whose realization requires the action of many other social and economic sectors in addition to the health sector." In its widest form the practice of medicine, that is to say the promotion and care of health, is concerned with this ideal.

ORGANIZATION OF HEALTH SERVICES

It is generally the goal of most countries to have their health services organized in such a way to ensure that individuals, families, and communities obtain the maximum benefit from current knowledge and technology available for the promotion, maintenance, and restoration of health. In order to play their part in this process,

governments and other agencies are faced with numerous tasks, including the following: (1) They must obtain as much information as is possible on the size, extent, and urgency of their needs; without accurate information, planning can be misdirected. (2) These needs must then be revised against the resources likely to be available in terms of money, manpower, and materials; developing countries may well require external aid to supplement their own resources. (3) Based on their assessments, countries then need to determine realistic objectives and draw up plans. (4) Finally, a process of evaluation needs to be built into the program; the lack of reliable information and accurate assessment can lead to confusion, waste, and inefficiency.

Health services of any nature reflect a number of interrelated characteristics, among which the most obvious, but not necessarily the most important from a national point of view, is the curative function; that is to say, caring for those already ill. Others include special services that deal with particular groups (such as children or pregnant women) and with specific needs such as nutrition or immunization; preventive services, the protection of the health both of individuals and of communities; health education; and, as mentioned above, the collection and analysis of information.

The Alma-Ata declaration

Levels of health care. In the curative domain there are various forms of medical practice. They may be thought of generally as forming a pyramidal structure, with three tiers representing increasing degrees of specialization and technical sophistication but catering to diminishing numbers of patients as they are filtered out of the system at a lower level. Only those patients who require special attention either for diagnosis or treatment should reach the second (advisory) or third (specialized treatment) tiers where the cost per item of service becomes increasingly higher. The first level represents primary health care, or first contact care, at which patients have their initial contact with the health-care system.

Primary health care is an integral part of a country's health maintenance system, of which it forms the largest and most important part. As described in the declaration of Alma-Ata, primary health care should be "based on practical, scientifically sound and socially acceptable methods and technology made universally accessible to individuals and families in the community through their full participation and at a cost that the community and country can afford to maintain at every stage of their development." Primary health care in the developed countries is usually the province of a medically qualified physician; in the developing countries first contact care is often provided by nonmedically qualified personnel.

The vast majority of patients can be fully dealt with at the primary level. Those who cannot are referred to the second tier (secondary health care, or the referral services) for the opinion of a consultant with specialized knowledge or for X-ray examinations and special tests. Secondary health care often requires the technology offered by a local or regional hospital. Increasingly, however, the radiological and laboratory services provided by hospitals are available directly to the family doctor, thus improving his service to patients and increasing its range. The third tier of health care, employing specialist services, is offered by institutions such as teaching hospitals and units devoted to the care of particular groups—women, children, patients with mental disorders, and so on. The dramatic differences in the cost of treatment at the various levels is a matter of particular importance in developing countries, where the cost of treatment for patients at the primary health-care level is usually only a small fraction of that at the third level; medical costs at any level in such countries, however, are usually borne by the government.

Ideally, provision of health care at all levels will be available to all patients; such health care may be said to be universal. The well-off, both in relatively wealthy industrialized countries and in the poorer developing world, may be able to get medical attention from sources they prefer and can pay for in the private sector. The vast majority of people in most countries, however, are dependent in various ways upon health services provided by the state, to which they may contribute comparatively little or, in the case of poor countries, nothing at all.

Costs of health care. The costs to national economies of providing health care are considerable and have been growing at a rapidly increasing rate, especially in countries such as the United States, Germany, and Sweden; the rise in Britain has been less rapid. This trend has been the cause of major concerns in both developed and developing countries. Some of this concern is based upon the lack of any consistent evidence to show that more spending on health care produces better health. There is a movement in developing countries to replace the type of organization of health-care services that evolved during European colonial times with some less expensive, and for them, more appropriate, health-care system.

In the industrialized world the growing cost of health services has caused both private and public health-care delivery systems to question current policies and to seek more economical methods of achieving their goals. Despite expenditures, health services are not always used effectively by those who need them, and results can vary widely from community to community. In Britain, for example, between 1951 and 1971 the death rate fell by 24 percent in the wealthier sections of the population but by only half that in the most underprivileged sections of

society. The achievement of good health is reliant upon more than just the quality of health care. Health entails such factors as good education, safe working conditions, a favourable environment, amenities in the home, well-integrated social services, and reasonable standards of living.

In the developing countries. The developing countries differ from one another culturally, socially, and economically, but what they have in common is a low average income per person, with large percentages of their populations living at or below the poverty level. Although most have a small elite class, living mainly in the cities, the largest part of their populations live in rural areas. Urban regions in developing and some developed countries in the mid- and late 20th century have developed pockets of slums, which are growing because of an influx of rural peoples. For lack of even the simplest measures, vast numbers of urban and rural poor die each year of preventable and curable diseases, often associated with poor hygiene and sanitation, impure water supplies, malnutrition, vitamin deficiencies, and chronic preventable infections. The effect of these and other deprivations is reflected by the finding that in the 1980s the life expectancy at birth for men and women was about one-third less in Africa than it was in Europe; similarly, infant mortality in Africa was about eight times greater than in Europe. The extension of primary health-care services is therefore a high priority in the developing countries.

The developing countries themselves, lacking the proper resources, have often been unable to generate or implement the plans necessary to provide required services at the village or urban poor level. It has, however, become clear that the system of health care that is appropriate for one country is often unsuitable for another. Research has established that effective health care is related to the special circumstances of the individual country, its people, culture, ideology, and economic and natural resources.

The rising costs of providing health care have influenced a trend, especially among the developing nations, to promote services that employ less highly trained primary health-care personnel who can be distributed more widely in order to reach the largest possible proportion of the community. The principal medical problems to be dealt with in the developing world include undernutrition, infection, gastrointestinal disorders, and respiratory complaints, which themselves may be the result of poverty, ignorance, and poor hygiene. For the most part, these are easy to identify and to treat. Furthermore, preventive measures are usually simple and cheap. Neither treatment nor prevention requires extensive professional training; in most cases they can be dealt with adequately by the "primary health worker," a term that includes all nonprofessional health personnel.

In the developed countries. Those concerned with providing health care in the developed countries face a different set of problems. The diseases so prevalent in the Third World have, for the most part, been eliminated or are readily treatable. Many of the adverse environmental conditions and public health hazards have been conquered. Social services of varying degrees of adequacy have been provided. Public funds can be called upon to support the cost of medical care, and there are a variety of private insurance plans available to the consumer. Nevertheless, the funds that a government can devote to health care are limited and the cost of modern medicine continues to increase, thus putting adequate medical services beyond the reach of many. Adding to the expense of modern medical practices is the increasing demand for greater funding of health education and preventive measures specifically directed toward the poor. (Ha.Sc.)

ADMINISTRATION OF PRIMARY HEALTH CARE

In many parts of the world, particularly in developing countries, people get their primary health care, or first-contact care, where available at all, from nonmedically qualified personnel; these cadres of medical auxiliaries are being trained in increasing numbers to meet overwhelming needs among rapidly growing populations. Even among the comparatively wealthy countries of the world, containing in all a much smaller percentage of the world's

Relating health-care systems to different cultures

Cost efficiency of health-care delivery

Medical auxiliary training

population, escalation in the costs of health services and in the cost of training a physician has precipitated some movement toward reappraisal of the role of the medical doctor in the delivery of first-contact care.

In advanced industrial countries, however, it is usually a trained physician who is called upon to provide the first-contact care. The patient seeking first-contact care can go either to a general practitioner or turn directly to a specialist. Which is the wisest choice has become a subject of some controversy. The general practitioner, however, is becoming rather rare in some developed countries. In countries where he does still exist, he is being increasingly observed as an obsolescent figure, because medicine covers an immense, rapidly changing, and complex field of which no physician can possibly master more than a small fraction. The very concept of the general practitioner, it is thus argued, may be absurd.

The obvious alternative to general practice is the direct access of a patient to a specialist. If a patient has problems with vision, he goes to an eye specialist, and if he has a pain in his chest (which he fears is due to his heart), he goes to a heart specialist. One objection to this plan is that the patient often cannot know which organ is responsible for his symptoms, and the most careful physician, after doing many investigations, may remain uncertain as to the cause. Breathlessness—a common symptom—may be due to heart disease, to lung disease, to anemia, or to emotional upset. Another common symptom is general malaise—feeling run-down or always tired; others are headache, chronic low backache, rheumatism, abdominal discomfort, poor appetite, and constipation. Some patients may also be overtly anxious or depressed. Among the most subtle medical skills is the ability to assess people with such symptoms and to distinguish between symptoms that are caused predominantly by emotional upset and those that are predominantly of bodily origin. A specialist may be capable of such a general assessment, but, often, with emphasis on his own subject, he fails at this point. The generalist with his broader training is often the better choice for a first diagnosis, with referral to a specialist as the next option.

It is often felt that there are also practical advantages for the patient in having his own doctor, who knows about his background, who has seen him through various illnesses, and who has often looked after his family as well. This personal physician, often a generalist, is in the best position to decide when the patient should be referred to a consultant.

The advantages of general practice and specialization are combined when the physician of first contact is a pediatrician. Although he sees only children and thus acquires a special knowledge of childhood maladies, he remains a generalist who looks at the whole patient. Another combination of general practice and specialization is represented by group practice, the members of which partially or fully specialize. One or more may be general practitioners, and one may be a surgeon, a second an obstetrician, a third a pediatrician, and a fourth an internist. In isolated communities group practice may be a satisfactory compromise, but in urban regions, where nearly everyone can be sent quickly to a hospital, the specialist surgeon working in a fully equipped hospital can usually provide better treatment than a general practitioner surgeon in a small clinic hospital.

MEDICAL PRACTICE IN DEVELOPED COUNTRIES

Britain. Before 1948, general practitioners in Britain settled where they could make a living. Patients fell into two main groups: weekly wage earners, who were compulsorily insured, were on a doctor's "panel" and were given free medical attention (for which the doctor was paid quarterly by the government); most of the remainder paid the doctor a fee for service at the time of the illness. In 1948 the National Health Service began operation. Under its provisions, everyone is entitled to free medical attention with a general practitioner with whom he is registered. Though general practitioners in the National Health Service are not debarred from also having private patients, these must be people who are not registered with

them under the National Health Service. Any physician is free to work as a general practitioner entirely independent of the National Health Service, though there are few who do so. Almost the entire population is registered with a National Health Service general practitioner, and the vast majority automatically sees this physician, or one of his partners, when they require medical attention. A few people, mostly wealthy, while registered with a National Health Service general practitioner, regularly see another physician privately; and a few may occasionally seek a private consultation because they are dissatisfied with their National Health Service physician.

A general practitioner under the National Health Service remains an independent contractor, paid by a capitation fee; that is, according to the number of people registered with him. He may work entirely from his own office, and he provides and pays his own receptionist, secretary, and other ancillary staff. Most general practitioners have one or more partners and work more and more in premises built for the purpose. Some of these structures are erected by the physicians themselves, but many are provided by the local authority, the physicians paying rent for using them. Health centres, in which groups of general practitioners work have become common.

In Britain only a small minority of general practitioners can admit patients to a hospital and look after them personally. Most of this minority are in country districts, where, before the days of the National Health Service, there were cottage hospitals run by general practitioners; many of these hospitals continued to function in a similar manner. All general practitioners use such hospital facilities as X-ray departments and laboratories, and many general practitioners work in hospitals in emergency rooms (casualty departments) or as clinical assistants to consultants, or specialists.

General practitioners are spread more evenly over the country than formerly, when there were many in the richer areas and few in the industrial towns. The maximum allowed list of National Health Service patients per doctor is 3,500; the average is about 2,500. Patients have free choice of the physician with whom they register, with the proviso that they cannot be accepted by one who already has a full list and that a physician can refuse to accept them (though such refusals are rare). In remote rural places there may be only one physician within a reasonable distance.

Until the mid-20th century it was not unusual for the doctor in Britain to visit patients in their own homes. A general practitioner might make 15 or 20 such house calls in a day, as well as seeing patients in his office or "surgery," often in the evenings. This enabled him to become a family doctor in fact as well as in name. In modern practice, however, a home visit is quite exceptional and is paid only to the severely disabled or seriously ill when other recourses are ruled out. All patients are normally required to go to the doctor.

It has also become unusual for a personal doctor to be available during weekends or holidays. His place may be taken by one of his partners in a group practice, a provision that is reasonably satisfactory. General practitioners, however, may now use one of several commercial deputizing services that employs young doctors to be on call. Although some of these young doctors may be well experienced, patients do not generally appreciate this kind of arrangement. (J.W.T./Ha.Sc.)

United States. Whereas in Britain the doctor of first contact is regularly a general practitioner, in the United States the nature of first-contact care is less consistent. General practice in the United States has been in a state of decline in the second half of the 20th century, especially in metropolitan areas. The general practitioner, however, is being replaced to some degree by the growing field of family practice. In 1969 family practice was recognized as a medical specialty after the American Academy of General Practice (now the American Academy of Family Physicians) and the American Medical Association created the American Board of General (now Family) Practice. Since that time the field has become one of the larger medical specialties in the United States. The family physi-

Home
visitations

Decline
of general
practice

cians were the first group of medical specialists in the United States for whom recertification was required.

There is no national health service, as such, in the United States. Most physicians in the country have traditionally been in some form of private practice, whether seeing patients in their own offices, clinics, medical centres, or another type of facility and regardless of the patients' income. Doctors are usually compensated by such state and federally supported agencies as Medicaid (for treating the poor) and Medicare (for treating the elderly); not all doctors, however, accept poor patients. There are also some state-supported clinics and hospitals where the poor and elderly may receive free or low-cost treatment, and some doctors devote a small percentage of their time to treatment of the indigent. Veterans may receive free treatment at Veterans Administration hospitals, and the federal government through its Indian Health Service provides medical services to American Indians and Alaskan natives, sometimes using trained auxiliaries for first-contact care.

In the rural United States first-contact care is likely to come from a generalist. The middle- and upper-income groups living in urban areas, however, have access to a larger number of primary medical care options. Children are often taken to pediatricians, who may oversee the child's health needs until adulthood. Adults frequently make their initial contact with an internist, whose field is mainly that of medical (as opposed to surgical) illnesses; the internist often becomes the family physician. Other adults choose to go directly to physicians with narrower specialties, including dermatologists, allergists, gynecologists, orthopedists, and ophthalmologists.

Patients in the United States may also choose to be treated by doctors of osteopathy. These doctors are fully qualified, but they make up only a small percentage of the country's physicians. They may also branch off into specialties, but general practice is much more common in their group than among M.D.'s.

It used to be more common in the United States for physicians providing primary care to work independently, providing their own equipment and paying their own ancillary staff. In smaller cities they mostly had full hospital privileges, but in larger cities these privileges were more likely to be restricted. Physicians, often sharing the same specialties, are increasingly entering into group associations, where the expenses of office space, staff, and equipment may be shared; such associations may work out of suites of offices, clinics, or medical centres. The increasing competition and risks of private practice have caused many physicians to join Health Maintenance Organizations (HMOs), which provide comprehensive medical care and hospital care on a prepaid basis. The cost savings to patients are considerable, but they must use only the HMO doctors and facilities. HMOs stress preventive medicine and out-patient treatment as opposed to hospitalization as a means of reducing costs, a policy that has caused an increased number of empty hospital beds in the United States.

While the number of doctors per 100,000 population in the United States has been steadily increasing, there has been a trend among physicians toward the use of trained medical personnel to handle some of the basic services normally performed by the doctor. So-called physician extender services are commonly divided into nurse practitioners and physician's assistants, both of whom provide similar ancillary services for the general practitioner or specialist. Such personnel do not replace the doctor. Almost all American physicians have systems for taking each other's calls when they become unavailable. House calls in the United States, as in Britain, have become exceedingly rare. (Ed.)

Russia. In Russia general practitioners are prevalent in the thinly populated rural areas. Pediatricians deal with children up to about age 15. Internists look after the medical ills of adults, and occupational physicians deal with the workers, sharing care with internists.

Teams of physicians with experience in varying specialties work from polyclinics or outpatient units, where many types of diseases are treated. Small towns usually have one polyclinic to serve all purposes. Large cities commonly

have separate polyclinics for children and adults, as well as clinics with specializations such as women's health care, mental illnesses, and sexually transmitted diseases. Polyclinics usually have X-ray apparatus and facilities for examination of tissue specimens, facilities associated with the departments of the district hospital. Beginning in the late 1970s was a trend toward the development of more large, multipurpose treatment centres, first-aid hospitals, and specialized medicine and health care centres.

Home visits have traditionally been common, and much of the physician's time is spent in performing routine checkups for preventive purposes. Some patients in sparsely populated rural areas may be seen first by feldshers (auxiliary health workers), nurses, or midwives who work under the supervision of a polyclinic or hospital physician. The feldsher was once a lower-grade physician in the army or peasant communities, but feldshers are now regarded as paramedical workers.

Japan. In Japan, with less rigid legal restriction of the sale of pharmaceuticals than in the West, there was formerly a strong tradition of self-medication and self-treatment. This was modified in 1961 by the institution of health insurance programs that covered a large proportion of the population; there was then a great increase in visits to the outpatient clinics of hospitals and to private clinics and individual physicians.

When Japan shifted from traditional Chinese medicine with the adoption of Western medical practices in the 1870s, Germany became the chief model. As a result of German influence and of their own traditions, Japanese physicians tended to prefer professorial status and scholarly research opportunities at the universities or positions in the national or prefectural hospitals to private practice. There were some pioneering physicians, however, who brought medical care to the ordinary people.

Physicians in Japan have tended to cluster in the urban areas. The Medical Service Law of 1963 was amended to empower the Ministry of Health and Welfare to control the planning and distribution of future public and non-profit medical facilities, partly to redress the urban-rural imbalance. Meanwhile, mobile services were expanded.

The influx of patients into hospitals and private clinics after the passage of the national health insurance acts of 1961 had, as one effect, a severe reduction in the amount of time available for any one patient. Perhaps in reaction to this situation, there has been a modest resurgence in the popularity of traditional Chinese medicine, with its leisurely interview, its dependence on herbal and other "natural" medicines, and its other traditional diagnostic and therapeutic practices. The rapid aging of the Japanese population as a result of the sharply decreasing death rate and birth rate has created an urgent need for expanded health care services for the elderly. There has also been an increasing need for centres to treat health problems resulting from environmental causes.

Other developed countries. On the continent of Europe there are great differences both within single countries and between countries in the kinds of first-contact medical care. General practice, while declining in Europe as elsewhere, is still rather common even in some large cities, as well as in remote country areas.

In The Netherlands, departments of general practice are administered by general practitioners in all the medical schools—an exceptional state of affairs—and general practice flourishes. In the larger cities of Denmark, general practice on an individual basis is usual and popular, because the physician works only during office hours. In addition, there is a duty doctor service for nights and weekends. In the cities of Sweden, primary care is given by specialists. In the remote regions of northern Sweden, district doctors act as general practitioners to patients spread over huge areas; the district doctors delegate much of their home visiting to nurses.

In France there are still general practitioners, but their number is declining. Many medical practitioners advertise themselves directly to the public as specialists in internal medicine, ophthalmologists, gynecologists, and other kinds of specialists. Even when patients have a general practitioner, they may still go directly to a specialist. At-

Health
Main-
tenance
Organiza-
tions

Poly-
clinics

Medical
practice in
Europe

tempts to stem the decline in general practice are being made by the development of group practice and of small rural hospitals equipped to deal with less serious illnesses, where general practitioners can look after their patients.

Although Israel has a high ratio of physicians to population, there is a shortage of general practitioners, and only in rural areas is general practice common. In the towns many people go directly to pediatricians, gynecologists, and other specialists, but there has been a reaction against this direct access to the specialist. More general practitioners have been trained, and the Israel Medical Association has recommended that no patient should be referred to a specialist except by the family physician or on instructions given by the family nurse. At Tel Aviv University there is a department of family medicine. In some newly developing areas, where the doctor shortage is greatest, there are medical centres at which all patients are initially interviewed by a nurse. The nurse may deal with many minor ailments, thus freeing the physician to treat the more seriously ill.

Nearly half the medical doctors in Australia are general practitioners—a far higher proportion than in most other advanced countries—though, as elsewhere, their numbers are declining. They tend to do far more for their patients than in Britain, many performing such operations as removal of the appendix, gallbladder, or uterus, operations that elsewhere would be carried out by a specialist surgeon. Group practices are common.

MEDICAL PRACTICE IN DEVELOPING COUNTRIES

China. Health services in China since the Cultural Revolution have been characterized by decentralization and dependence on personnel chosen locally and trained for short periods. Emphasis is given to selfless motivation, self-reliance, and to the involvement of everyone in the community. Campaigns stressing the importance of preventive measures and their implementation have served to create new social attitudes as well as to break down divisions between different categories of health workers. Health care is regarded as a local matter that should not require the intervention of any higher authority; it is based upon a highly organized and well-disciplined system that is egalitarian rather than hierarchical, as in Western societies, and which is well suited to the rural areas where about two-thirds of the population live. In the large and crowded cities an important constituent of the health-care system is the residents' committees, each for a population of 1,000 to 5,000 people. Care is provided by part-time personnel with periodic visits by a doctor. A number of residents' committees are grouped together into neighbourhoods of some 50,000 people where there are clinics and general hospitals staffed by doctors as well as health auxiliaries trained in both traditional and Westernized medicine. Specialized care is provided at the district level (over 100,000 people), in district hospitals and in epidemic and preventive medicine centres. In many rural districts people's communes have organized cooperative medical services that provide primary care for a small annual fee.

Throughout China the value of traditional medicine is stressed, especially in the rural areas. All medical schools are encouraged to teach traditional medicine as part of their curriculum, and efforts are made to link colleges of Chinese medicine with Western-type medical schools. Medical education is of shorter duration than it is in Europe, and there is greater emphasis on practical work. Students spend part of their time away from the medical school working in factories or in communes; they are encouraged to question what they are taught and to participate in the educational process at all stages. One well-known form of traditional medicine is acupuncture, which is used as a therapeutic and pain-relieving technique; requiring the insertion of brass-handled needles at various points on the body, acupuncture has become quite prominent as a form of anesthesia.

The vast number of nonmedically qualified health staff, upon whom the health-care system greatly depends, includes both full-time and part-time workers. The latter include so-called barefoot doctors, who work mainly in rural areas, worker doctors in factories, and medical work-

ers in residential communities. None of these groups is medically qualified. They have had only a three-month period of formal training, part of which is done in a hospital, fairly evenly divided between theoretical and practical work. This is followed by a varying period of on-the-job experience under supervision.

India. Āyurvedic medicine is an example of a well-organized system of traditional health care, both preventive and curative, that is widely practiced in parts of Asia. Āyurvedic medicine has a long tradition behind it, having originated in India perhaps as long as 3,000 years ago. It is still a favoured form of health care in large parts of the Eastern world, especially in India, where a large percentage of the population use this system exclusively or combined with modern medicine. The Indian Medical Council was set up in 1971 by the Indian government to establish maintenance of standards for undergraduate and postgraduate education. It establishes suitable qualifications in Indian medicine and recognizes various forms of traditional practice including Āyurvedic, Unani, and Siddha. Projects have been undertaken to integrate the indigenous Indian and Western forms of medicine. Most Āyurvedic practitioners work in rural areas, providing health care to at least 500,000,000 people in India alone. They therefore represent a major force for primary health care, and their training and deployment are important to the government of India.

Like scientific medicine, Āyurvedic medicine has both preventive and curative aspects. The preventive component emphasizes the need for a strict code of personal and social hygiene, the details of which depend upon individual, climatic, and environmental needs. Bodily exercises, the use of herbal preparations, and Yoga form a part of the remedial measures. The curative aspects of Āyurvedic medicine involves the use of herbal medicines, external preparations, physiotherapy, and diet. It is a principle of Āyurvedic medicine that the preventive and therapeutic measures be adapted to the personal requirements of each patient.

Other developing countries. A main goal of the World Health Organization (WHO), as expressed in the Alma-Ata Declaration of 1978, is to provide to all the citizens of the world a level of health that will allow them to lead socially and economically productive lives by the year 2000. By the late 1980s, however, vast disparities in health care still existed between the rich and poor countries of the world. In developing countries such as Ethiopia, Guinea, Mali, and Mozambique, for instance, governments in the late 1980s spent less than \$5 per person per year on public health, while in most western European countries several hundred dollars per year was spent on each person. The disproportion of the number of physicians available between developing and developed countries is similarly wide.

Along with the shortage of physicians, there is a shortage of everything else needed to provide medical care—of equipment, drugs, and suitable buildings, and of nurses, technicians, and all other grades of staff, whose presence is taken for granted in the affluent societies. Yet there are greater percentages of sick in the poor countries than in the rich countries. In the poor countries a high proportion of people are young, and all are liable to many infections, including tuberculosis, syphilis, typhoid, and cholera (which, with the possible exception of syphilis, are now rare in the rich countries), and also malaria, yaws, worm infestations, and many other conditions occurring primarily in the warmer climates. Nearly all of these infections respond to the antibiotics and other drugs that have been discovered since the 1920s. There is also much malnutrition and anemia, which can be cured if funding is available. There is a prevalence of disorders remediable by surgery. Preventive medicine can ensure clean water supplies, destroy insects that carry infections, teach hygiene, and show how to make the best use of resources.

In most poor countries there are a few people, usually living in the cities, who can afford to pay for medical care, and in a free market system the physicians tend to go where they can make the best living; this situation causes the doctor-patient ratio to be much higher in the towns

Āyurvedic
medicine

Residents'
committees

Acu-
puncture

Medical
needs in
developing
countries

than in country districts. A physician in Bombay or in Rio de Janeiro, for example, may have equipment as lavish as that of a physician in the United States and can earn an excellent income. The poor, however, both in the cities and in the country, can get medical attention only if it is paid for by the state, by some supranational body, or by a mission or other charitable organization. Moreover, the quality of the care they receive is often poor, and in remote regions it may be lacking altogether. In practice, hospitals run by a mission may cooperate closely with state-run health centres.

Because physicians are scarce, their skills must be used to best advantage, and much of the work normally done by physicians in the rich countries has to be delegated to auxiliaries or nurses, who have to diagnose the common conditions, give treatment, take blood samples, help with operations, supply simple posters containing health advice, and carry out other tasks. In such places the doctor has time only to perform major operations and deal with the more difficult medical problems. People are treated as far as possible on an outpatient basis from health centres housed in simple buildings; few can travel except on foot, and, if they are more than a few miles from a health centre, they tend not to go there. Health centres also may be used for health education.

Health
care in
Tanzania

Although primary health-care service differs from country to country, that developed in Tanzania is representative of many that have been devised in largely rural developing countries. The most important feature of the Tanzanian rural health service is the rural health centre, which, with its related dispensaries, is intended to provide comprehensive health services for the community. The staff is headed by the assistant medical officer and the medical assistant. The assistant medical officer has at least four years of experience, which is then followed by further training for 18 months. He is not a doctor but serves to bridge the gap between medical assistant and physician. The medical assistant has three years of general medical education. The work of the rural health centres and dispensaries is mainly of three kinds: diagnosis and treatment, maternal and child health, and environmental health. The main categories of primary health workers also include medical aids, maternal and child health aids, and health auxiliaries. Nurses and midwives form another category of worker. In the villages there are village health posts staffed by village medical helpers working under supervision from the rural health centre.

Primitive
medicine

In some primitive elements of the societies of developing countries, and of some developed countries, there exists the belief that illness comes from the displeasure of ancestral gods and evil spirits, from the malign influence of evilly disposed persons, or from natural phenomena that can neither be forecast nor controlled. To deal with such causes there are many varieties of indigenous healers who practice elaborate rituals on behalf of both the physically ill and the mentally afflicted. If it is understood that such beliefs, and other forms of shamanism, may provide a basis upon which health care can be based, then primary health care may be said to exist almost everywhere. It is not only easily available but also readily acceptable, and often preferred, to more rational methods of diagnosis and treatment. Although such methods may sometimes be harmful, they may often be effective, especially where the cause is psychosomatic. Other patients, however, may suffer from a disease for which there is a cure in modern medicine.

In order to improve the coverage of primary health-care services and to spread more widely some of the benefits of Western medicine, attempts have sometimes been made to find a means of cooperation, or even integration, between traditional and modern medicine (see above *India*). In Africa, for example, some such attempts are officially sponsored by ministries of health, state governments, universities, and the like, and they have the approval of WHO, which often takes the lead in this activity. In view, however, of the historical relationships between these two systems of medicine, their different basic concepts, and the fact that their methods cannot readily be combined, successful merging has been limited.

ALTERNATIVE OR COMPLEMENTARY MEDICINE

Persons dissatisfied with the methods of modern medicine or with its results sometimes seek help from those professing expertise in other, less conventional, and sometimes controversial, forms of health care. Such practitioners are not medically qualified unless they are combining such treatments with a regular (allopathic) practice, which includes osteopathy. In many countries the use of some forms, such as chiropractic, requires licensing and a degree from an approved college. The treatments afforded in these various practices are not always subjected to objective assessment, yet they provide services that are alternative, and sometimes complementary, to conventional practice. This group includes practitioners of homeopathy, naturopathy, acupuncture, hypnotism, and various meditative and quasi-religious forms. Numerous persons also seek out some form of faith healing to cure their ills, sometimes as a means of last resort. Religions commonly include some advents of miraculous curing within their scriptures. The belief in such curative powers has been in part responsible for the increasing popularity of the television, or "electronic," preacher in the United States, a phenomenon that involves millions of viewers. Millions of others annually visit religious shrines, such as the one at Lourdes in France, with the hope of being miraculously healed.

Faith
healing

SPECIAL PRACTICES AND FIELDS OF MEDICINE

Specialties in medicine. At the beginning of World War II it was possible to recognize a number of major medical specialties, including internal medicine, obstetrics and gynecology, pediatrics, pathology, anesthesiology, ophthalmology, surgery, orthopedic surgery, plastic surgery, psychiatry and neurology, radiology, and urology. Hematology was also an important field of study, and microbiology and biochemistry were important medically allied specialties. Since World War II, however, there has been an almost explosive increase of knowledge in the medical sciences as well as enormous advances in technology as applicable to medicine. These developments have led to more and more specialization. The knowledge of pathology has been greatly extended, mainly by the use of the electron microscope; similarly microbiology, which includes bacteriology, expanded with the growth of such other subfields as virology (the study of viruses) and mycology (the study of yeasts and fungi in medicine). Biochemistry, sometimes called clinical chemistry or chemical pathology, has contributed to the knowledge of disease, especially in the field of genetics where genetic engineering has become a key to curing some of the most difficult diseases. Hematology also expanded after World War II with the development of electron microscopy. Contributions to medicine have come from such fields as psychology and sociology especially in such areas as mental disorders and mental handicaps. Clinical pharmacology has led to the development of more effective drugs and to the identification of adverse reactions. More recently established medical specialties are those of preventive medicine, physical medicine and rehabilitation, family practice, and nuclear medicine. In the United States every medical specialist must be certified by a board composed of members of the specialty in which certification is sought. Some type of peer certification is required in most countries.

Certifica-
tion

Expansion of knowledge both in depth and in range has encouraged the development of new forms of treatment that require high degrees of specialization, such as organ transplantation and exchange transfusion; the field of anesthesiology has grown increasingly complex as equipment and anesthetics have improved. New technologies have introduced microsurgery, laser beam surgery, and lens implantation (for cataract patients), all requiring the specialist's skill. Precision in diagnosis has markedly improved; advances in radiology, the use of ultrasound, computerized axial tomography (CAT scan), and nuclear magnetic resonance imaging are examples of the extension of technology requiring expertise in the field of medicine.

To provide more efficient service it is not uncommon for a specialist surgeon and a specialist physician to form a team working together in the field of, for example, heart

disease. An advantage of this arrangement is that they can attract a highly trained group of nurses, technologists, operating room technicians, and so on, thus greatly improving the efficiency of the service to the patient. Such specialization is expensive, however, and has required an increasingly large proportion of the health budget of institutions, a situation that eventually has its financial effect on the individual citizen. The question therefore arises as to their cost-effectiveness. Governments of developing countries have usually found, for instance, that it is more cost-efficient to provide more people with basic care.

Advantage
of medical
school
association

Teaching. Physicians in developed countries frequently prefer posts in hospitals with medical schools. Newly qualified physicians want to work there because doing so will aid their future careers, though the actual experience may be wider and better in a hospital without a medical school. Senior physicians seek careers in hospitals with medical schools because consultant, specialist, or professorial posts there usually carry a high degree of prestige. When the posts are salaried, the salaries are sometimes, but not always, higher than in a nonteaching hospital. Usually a consultant who works in private practice earns more when on the staff of a medical school.

In many medical schools there are clinical professors in each of the major specialties—such as surgery, internal medicine, obstetrics and gynecology, and psychiatry—and often of the smaller specialties as well. There are also professors of pathology, radiology, and radiotherapy. Whether professors or not, all doctors in teaching hospitals have the two functions of caring for the sick and educating students. They give lectures and seminars and are accompanied by students on ward rounds.

Industrial medicine. The Industrial Revolution greatly changed, and as a rule worsened, the health hazards caused by industry, while the numbers at risk vastly increased. In Britain the first small beginnings of efforts to ameliorate the lot of the workers in factories and mines began in 1802 with the passing of the first factory act, the Health and Morals of Apprentices Act. The factory act of 1838, however, was the first truly effective measure in the industrial field. It forbade night work for children and restricted their work hours to 12 per day. Children under 13 were required to attend school. A factory inspectorate was established, the inspectors being given powers of entry into factories and power of prosecution of recalcitrant owners. Thereafter there was a succession of acts with detailed regulations for safety and health in all industries. Industrial diseases were made notifiable, and those who developed any prescribed industrial disease were entitled to benefits.

Industrial
physician's
function

The situation is similar in other developed countries. Physicians are bound by legal restrictions and must report industrial diseases. The industrial physician's most important function, however, is to prevent industrial diseases. Many of the measures to this end have become standard practice, but, especially in industries working with new substances, the physician should determine if workers are being damaged and suggest preventive measures. The industrial physician may advise management about industrial hygiene and the need for safety devices and protective clothing and may become involved in building design. The physician or health worker may also inform the worker of occupational health hazards.

Modern factories usually have arrangements for giving first aid in case of accidents. Depending upon the size of the plant, the facilities may range from a simple first-aid station to a large suite of lavishly equipped rooms and may include a staff of qualified nurses and physiotherapists and one or perhaps more full-time physicians.

Periodic medical examination. Physicians in industry carry out medical examinations, especially on new employees and on those returning to work after sickness or injury. In addition, those liable to health hazards may be examined regularly in the hope of detecting evidence of incipient damage. In some organizations every employee may be offered a regular medical examination.

The industrial and the personal physician. When a worker also has a personal physician, there may be doubt, in some cases, as to which physician bears the main responsibility for his health. When someone has an accident

or becomes acutely ill at work, the first aid is given or directed by the industrial physician. Subsequent treatment may be given either at the clinic at work or by the personal physician. Because of labour-management difficulties, workers sometimes tend not to trust the diagnosis of the management-hired physician.

Industrial health services. During the epoch of the Soviet Union and the Soviet bloc, industrial health service generally developed more fully in those countries than in the capitalist countries. At the larger industrial establishments in the Soviet Union, polyclinics were created to provide both occupational and general care for workers and their families. Occupational physicians were responsible for preventing occupational diseases and injuries, health screening, immunization, and health education.

In the capitalist countries, on the other hand, no fixed pattern of industrial health service has emerged. Legislation impinges upon health in various ways, including the provision of safety measures, the restriction of pollution, and the enforcement of minimum standards of lighting, ventilation, and space per person. In most of these countries there is found an infinite variety of schemes financed and run by individual firms or, equally, by huge industries. Labour unions have also done much to enforce health codes within their respective industries. In the developing countries there has been generally little advance in industrial medicine.

Family health care. In many societies special facilities are provided for the health care of pregnant women, mothers, and their young children. The health care needs of these three groups are generally recognized to be so closely related as to require a highly integrated service that includes prenatal care, the birth of the baby, the postnatal period, and the needs of the infant. Such a continuum should be followed by a service attentive to the needs of young children and then by a school health service. Family clinics are common in countries that have state-sponsored health services, such as those in the United Kingdom and elsewhere in Europe. Family health care in some developed countries, such as the United States, is provided for low-income groups by state-subsidized facilities, but other groups defer to private physicians or privately run clinics.

The
prenatal
clinic

Prenatal clinics provide a number of elements. There is, first, the care of the pregnant woman, especially if she is in a vulnerable group likely to develop some complication during the last few weeks of pregnancy and subsequent delivery. Many potential hazards, such as diabetes and high blood pressure, can be identified and measures taken to minimize their effects. In developing countries pregnant women are especially susceptible to many kinds of disorders, particularly infections such as malaria. Local conditions determine what special precautions should be taken to ensure a healthy child. Most pregnant women, in their concern to have a healthy child, are receptive to simple health education. The prenatal clinic provides an excellent opportunity to teach the mother how to look after herself during pregnancy, what to expect at delivery, and how to care for her baby. If the clinic is attended regularly, the woman's record will be available to the staff that will later supervise the delivery of the baby; this is particularly important for someone who has been determined to be at risk. The same clinical unit should be responsible for prenatal, natal, and postnatal care as well as for the care of the newborn infants.

Most pregnant women can be safely delivered in simple circumstances without an elaborately trained staff or sophisticated technical facilities, provided that these can be called upon in emergencies. In developed countries it was customary in premodern times for the delivery to take place in the woman's home supervised by a qualified midwife or by the family doctor. By the mid-20th century women, especially in urban areas, usually preferred to have their babies in a hospital, either in a general hospital or in a more specialized maternity hospital. In many developing countries traditional birth attendants supervise the delivery. They are women, for the most part without formal training, who have acquired skill by working with others and from their own experience. Normally they belong to the local community where they have the confidence of

the family, where they are content to live and serve, and where their services are of great value. In many developing countries the better training of birth attendants has a high priority. In developed Western countries there has been a trend toward delivery by natural childbirth, including delivery in a hospital without anesthesia, and home delivery.

Postnatal health care

Postnatal care services are designed to supervise the return to normal of the mother. They are usually given by the staff of the same unit that was responsible for the delivery. Important considerations are the matter of breast- or artificial feeding and the care of the infant. Today the prospects for survival of babies born prematurely or after a difficult and complicated labour, as well as for neonates (recently born babies) with some physical abnormality, are vastly improved. This is due to technical advances, including those that can determine defects in the prenatal stage, as well as to the growth of neonatology as a specialty. A vital part of the family health-care service is the child welfare clinic, which undertakes the care of the newborn. The first step is the thorough physical examination of the child on one or more occasions to determine whether or not it is normal both physically and, if possible, mentally. Later periodic examinations serve to decide if the infant is growing satisfactorily. Arrangements can be made for the child to be protected from major hazards by, for example, immunization and dietary supplements. Any intercurrent condition, such as a chest infection or skin disorder, can be detected early and treated. Throughout the whole of this period mother and child are together, and particular attention is paid to the education of the mother for the care of the child.

A part of the health service available to children in the developed countries is that devoted to child guidance. This provides psychiatric guidance to maladjusted children usually through the cooperative work of a child psychiatrist, educational psychologist, and schoolteacher.

Increasing need for elderly care

Geriatrics. Since the mid-20th century a change has occurred in the population structure in developed countries. The proportion of elderly people has been increasing. Since 1983, however, in most European countries the population growth of that group has leveled off, although it is expected to continue to grow more rapidly than the rest of the population in most countries through the first third of the 21st century. In the late 20th century Japan had the fastest growing elderly population.

Geriatrics, the health care of the elderly, is therefore a considerable burden on health services. In the United Kingdom about one-third of all hospital beds are occupied by patients over 65; half of these are psychiatric patients. The physician's time is being spent more and more with the elderly, and since statistics show that women live longer than men, geriatric practice is becoming increasingly concerned with the treatment of women. Elderly people often have more than one disorder, many of which are chronic and incurable, and they need more attention from health-care services. In the United States there has been some movement toward making geriatrics a medical specialty, but it has not generally been recognized.

Support services for the elderly provided by private or state-subsidized sources include domestic help, delivery of meals, day-care centres, elderly residential homes or nursing homes, and hospital beds either in general medical wards or in specialized geriatric units. The degree of accessibility of these services is uneven from country to country and within countries. In the United States, for instance, although there are some federal programs, each state has its own elderly programs, which vary widely. However, as the elderly become an increasingly larger part of the population their voting rights are providing increased leverage for obtaining more federal and state benefits. The general practitioner or family physician working with visiting health and social workers and in conjunction with the patient's family often form a working team for elderly care.

In the developing world, countries are largely spared such geriatric problems, but not necessarily for positive reasons. A principal cause, for instance, is that people do not live so long. Another major reason is that in the extended family concept, still prevalent among developing countries, most of the caretaking needs of the elderly are provided by the family.

Public health practice. The physician working in the field of public health is mainly concerned with the environmental causes of ill health and in their prevention. Bad drainage, polluted water and atmosphere, noise and smells, infected food, bad housing, and poverty in general are all his special concern. Perhaps the most descriptive title he can be given is that of community physician. In Britain he has been customarily known as the medical officer of health and, in the United States, as the health officer.

The spectacular improvement in the expectation of life in the affluent countries has been due far more to public health measures than to curative medicine. These public health measures began operation largely in the 19th century. At the beginning of that century, drainage and water supply systems were all more or less primitive; nearly all the cities of that time had poorer water and drainage systems than Rome had possessed 1,800 years previously. Infected water supplies caused outbreaks of typhoid, cholera, and other waterborne infections. By the end of the century, at least in the larger cities, water supplies were usually safe. Food-borne infections were also drastically reduced by the enforcement of laws concerned with the preparation, storage, and distribution of food. Insect-borne infections, such as malaria and yellow fever, which were common in tropical and semitropical climates, were eliminated by the destruction of the responsible insects. Fundamental to this improvement in health has been the diminution of poverty, for most public health measures are expensive. The peoples of the developing countries fall sick and sometimes die from infections that are virtually unknown in affluent countries.

Britain. Public health services in Britain are organized locally under the National Health Service. The medical officer of health is employed by the local council and is the adviser in health matters. The larger councils employ a number of mostly full-time medical officers; in some rural areas, a general practitioner may be employed part-time as medical officer of health.

The medical officer has various statutory powers conferred by acts of Parliament, regulations and orders, such as food and drugs acts, milk and dairies regulations, and factories acts. He supervises the work of sanitary inspectors in the control of health nuisances. The compulsorily notifiable infectious diseases are reported to him, and he takes appropriate action. Other concerns of the medical officer include those involved with the work of the district nurse, who carries out nursing duties in the home, and the health visitor, who gives advice on health matters, especially to the mothers of small babies. He has other duties in connection with infant welfare clinics, crèches, day and residential nurseries, the examination of schoolchildren, child guidance clinics, foster homes, factories, problem families, and the care of the aged and the handicapped.

United States. Federal, state, county, and city governments all have public health functions. Under the U.S. Department of Health and Human Services is the Public Health Service, headed by an assistant secretary for health and the surgeon general. State health departments are headed by a commissioner of health, usually a physician, who is often in the governor's cabinet. He usually has a board of health that adopts health regulations and holds hearings on their alleged violations. A state's public health code is the foundation on which all county and city health regulations must be based. A city health department may be independent of its surrounding county health department, or there may be a combined city-county health department. The physicians of the local health departments are usually called health officers, though occasionally people with this title are not physicians. The larger departments may have a public health director, a district health director, or a regional health director.

The minimal complement of a local health department is a health officer, a public health nurse, a sanitation expert, and a clerk who is also a registrar of vital statistics. There may also be sanitation personnel, nutritionists, social workers, laboratory technicians, and others.

Japan. Japan's Ministry of Health and Welfare directs public health programs at the national level, maintaining close coordination among the fields of preventive

Public health measures of the 19th century

Medical officer of health

medicine, medical care, and welfare and health insurance. The departments of health of the prefectures and of the largest municipalities operate health centres. The integrated community health programs of the centres encompass maternal and child health, communicable-disease control, health education, family planning, health statistics, food inspection, and environmental sanitation. Private physicians, through their local medical associations, help to formulate and execute particular public health programs needed by their localities.

Numerous laws are administered through the ministry's bureaus and agencies, which range from public health, environmental sanitation, and medical affairs to the children and families bureau. The various categories of institutions run by the ministry, in addition to the national hospitals, include research centres for cancer and leprosy, homes for the blind, rehabilitation centres for the physically handicapped, and port quarantine services.

Former Soviet Union. In the aftermath of the dissolution of the Soviet Union, responsibility for public health fell to the governments of the successor countries.

The public health services for the U.S.S.R. as a whole were directed by the Ministry of Health. The ministry, through the 15 union republic ministries of health, directed all medical institutions within its competence as well as the public health authorities and services throughout the country.

The administration was centralized, with little local autonomy. Each of the 15 republics had its own ministry of health, which was responsible for carrying out the plans and decisions established by the U.S.S.R. Ministry of Health. Each republic was divided into *oblasti*, or provinces, which had departments of health directly responsible to the republic ministry of health. Each *oblast*, in turn, had *rayony* (municipalities), which have their own health departments accountable to the *oblast* health department. Finally, each *rayon* was subdivided into smaller *uchastoki* (districts).

In most rural *rayony* the responsibility for public health lay with the chief physician, who was also medical director of the central *rayon* hospital. This system ensured unity of public health administration and implementation of the principle of planned development. Other health personnel included nurses, feldshers, and midwives.

For more information on the history, organization, and progress of public health, see below.

Military practice. The medical services of armies, navies, and air forces are geared to war. During campaigns the first requirement is the prevention of sickness. In all wars before the 20th century, many more combatants died of disease than of wounds. And even in World War II and wars thereafter, although few died of disease, vast numbers became casualties from disease.

The main means of preventing sickness are the provision of adequate food and pure water, thus eliminating starvation, avitaminosis, and dysentery and other bowel infections, which used to be particular scourges of armies; the provision of proper clothing and other means of protection from the weather; the elimination from the service of those likely to fall sick; the use of vaccination and suppressive drugs to prevent various infections, such as typhoid and malaria; and education in hygiene and in the prevention of sexually transmitted diseases, a particular problem in the services. In addition, the maintenance of high morale has a striking effect on casualty rates, for, when morale is poor, soldiers are likely to suffer psychiatric breakdowns, and malingering is more prevalent.

The medical branch may provide advice about disease prevention, but the actual execution of this advice is through the ordinary chains of command. It is the duty of the military, not of the medical, officer to ensure that the troops obey orders not to drink infected water and to take tablets to suppress malaria.

Army medical organization. The medical doctor of first contact to the soldier in the armies of developed countries is usually an officer in the medical corps. In peacetime the doctor sees the sick and has functions similar to those of the general practitioner, prescribing drugs and dressings, and there may be a sick bay where slightly sick soldiers can

remain for a few days. The doctor is usually assisted by trained nurses and corpsmen. If a further medical opinion is required, the patient can be referred to a specialist at a military or civilian hospital.

In a war zone, medical officers have an aid post where, with the help of corpsmen, they apply first aid to the walking wounded and to the more seriously wounded who are brought in. The casualties are evacuated as quickly as possible by field ambulances or helicopters. At a company station, medical officers and medical corpsmen may provide further treatment before patients are evacuated to the main dressing station at the field ambulance headquarters, where a surgeon may perform emergency operations. Thereafter, evacuation may be to casualty clearing stations, to advanced hospitals, or to base hospitals. Air evacuation is widely used.

In peacetime most of the intermediate medical units exist only in skeleton form; the active units are at the battalion and hospital level. When physicians join the medical corps, they may join with specialist qualifications, or they may obtain such qualifications while in the army. A feature of army medicine is promotion to administrative positions. The commanding officer of a hospital and the medical officer at headquarters may have no contacts with actual patients.

Although medical officers in peacetime have some choice of the kind of work they will do, they are in a chain of command and are subject to military discipline. When dealing with patients, however, they are in a special position; they cannot be ordered by a superior officer to give some treatment or take other action that they believe is wrong. Medical officers also do not bear or use arms unless their patients are being attacked.

Naval and air force medicine. Naval medical services are run on lines similar to those of the army. Junior medical officers are attached to ships or to shore stations and deal with most cases of sickness in their units. When at sea, medical officers have an exceptional degree of responsibility in that they work alone, unless they are on a very large ship. In peacetime, only the larger ships carry a medical officer; in wartime, destroyers and other small craft may also carry medical officers. Serious cases go to either a shore-based hospital or a hospital ship.

Flying has many medical repercussions. Cold, lack of oxygen, and changes of direction at high speed all have important effects on bodily and mental functions. Armies and air forces may share the same medical services.

A developing field is aerospace medicine. This involves medical problems that were not experienced before spaceflight, for the main reason that humans in space are not under the influence of gravity, a condition that has profound physiological effects.

CLINICAL RESEARCH

The remarkable developments in medicine that have been brought about in the 20th century, especially since World War II, have been based on research either in the basic sciences related to medicine or in the clinical field. Advances in the use of radiation, nuclear energy, and space research have played an important part in this progress. Some laypersons often think of research as taking place only in sophisticated laboratories or highly specialized institutions where work is devoted to scientific advances that may or may not be applicable to medical practice. This notion, however, ignores the clinical research that takes place on a day-to-day basis in hospitals and doctors' offices.

Historical notes. Although the most spectacular changes in the medical scene during the 20th century, and the most widely heralded, have been the development of potent drugs and elaborate operations, another striking change has been the abandonment of most of the remedies of the past. In the mid-19th century, persons ill with numerous maladies were starved (partially or completely), bled, purged, cupped (by applying a tight-fitting vessel filled with steam to some part and then cooling the vessel), and rested, perhaps for months or even years. Much more recently they were prescribed various restricted diets and were routinely kept in bed for weeks after abdominal operations, for many weeks or months when their hearts were

Dealing with the wounded

Medical organization in the navy

Military prevention of sickness

Abandonment of remedies of the past

thought to be affected, and for many months or years with tuberculosis. The abandonment of these measures may not be thought of as involving research, but the physician who first encouraged persons who had peptic ulcers to eat normally (rather than to live on the customary bland foods) and the physician who first got his patients out of bed a week or two after they had had minor coronary thrombosis (rather than insisting on a minimum of six weeks of strict bed rest) were as much doing research as is the physician who first tries out a new drug on a patient. This research, by observing what happens when remedies are abandoned, has been of inestimable value, and the need for it has not passed.

Clinical observation. Much of the investigative clinical field work undertaken in the present day requires only relatively simple laboratory facilities because it is observational rather than experimental in character. A feature of much contemporary medical research is that it requires the collaboration of a number of persons, perhaps not all of them doctors. Despite the advancing technology, there is much to be learned simply from the observation and analysis of the natural history of disease processes as they begin to affect patients, pursue their course, and end, either in their resolution or by the death of the patient. Such studies may be suitably undertaken by physicians working in their offices who are in a better position than doctors working only in hospitals to observe the whole course of an illness. Disease rarely begins in a hospital and usually does not end there. It is notable, however, that observational research is subject to many limitations and pitfalls of interpretation, even when it is carefully planned and meticulously carried out.

Drug research. The administration of any medicament, especially a new drug, to a patient is fundamentally an experiment: so is a surgical operation, particularly if it involves a modification to an established technique or a completely new procedure. Concern for the patient, careful observation, accurate recording, and a detached mind are the keys to this kind of investigation, as indeed to all forms of clinical study. Because patients are individuals reacting to a situation in their own different ways, the data obtained in groups of patients may well require statistical analysis for their evaluation and validation.

One of the striking characteristics in the medical field in the 20th century has been the development of new drugs, usually by pharmaceutical companies. Until the end of the 19th century, the discovery of new drugs was largely a matter of chance. It was in that period that Paul Ehrlich, the German scientist, began to lay down the principles for modern pharmaceutical research that made possible the development of a vast array of safe and effective drugs. Such benefits, however, bring with them their own disadvantages: it is estimated that as many as 30 percent of patients in, or admitted to, hospitals suffer from the adverse effect of drugs prescribed by a physician for their treatment (iatrogenic disease). Sometimes it is extremely difficult to determine whether a drug has been responsible for some disorder. An example of the difficulty is provided by the thalidomide disaster between 1959 and 1962. Only after numerous deformed babies had been born throughout the world did it become clear that thalidomide taken by the mother as a sedative had been responsible.

In hospitals where clinical research is carried out, ethical committees often consider each research project. If the committee believes that the risks are not justified, the project is rejected.

After a potentially useful chemical compound has been identified in the laboratory, it is extensively tested in animals, usually for a period of months or even years. Few drugs make it beyond this point. If the tests are satisfactory, the decision may be made for testing the drug in humans. It is this activity that forms the basis of much clinical research. In most countries the first step is the study of its effects in a small number of health volunteers. The response, effect on metabolism, and possible toxicity are carefully monitored and have to be completely satisfactory before the drug can be passed for further studies, namely with patients who have the disorder for which the drug is to be used. Tests are administered at first to a limited

number of these patients to determine effectiveness, proper dosage, and possible adverse reactions. These searching studies are scrupulously controlled under stringent conditions. Larger groups of patients are subsequently involved to gain a wider sampling of the information. Finally, a full-scale clinical trial is set up. If the regulatory authority is satisfied about the drug's quality, safety, and efficacy, it receives a license to be produced. As the drug becomes more widely used, it eventually finds its proper place in therapeutic practice, a process that may take years.

An important step forward in clinical research was taken in the mid-20th century with the development of the controlled clinical trial. This sets out to compare two groups of patients, one of which has had some form of treatment that the other group has not. The testing of a new drug is a case in point: one group receives the drug, the other a product identical in appearance, but which is known to be inert—a so-called placebo. At the end of the trial, the results of which can be assessed in various ways, it can be determined whether or not the drug is effective and safe. By the same technique two treatments can be compared, for example a new drug against a more familiar one. Because individuals differ physiologically and psychologically, the allocation of patients between the two groups must be made in a random fashion; some method independent of human choice must be used so that such differences are distributed equally between the two groups.

In order to reduce bias and make the trial as objective as possible the double-blind technique is sometimes used. In this procedure, neither the doctor nor the patients know which of two treatments is being given. Despite such precautions the results of such trials can be prejudiced, so that rigorous statistical analysis is required. It is obvious that many ethical, not to say legal, considerations arise, and it is essential that all patients have given their informed consent to be included. Difficulties arise when patients are unconscious, mentally confused, or otherwise unable to give their informed consent. Children present a special difficulty because not all laws agree that parents can legally commit a child to an experimental procedure. Trials, and indeed all forms of clinical research that involve patients, must often be submitted to a committee set up locally to scrutinize each proposal.

Surgery. In drug research the essential steps are taken by the chemists who synthesize or isolate new drugs in the laboratory; clinicians play only a subsidiary part. In developing new surgical operations clinicians play a more important role, though laboratory scientists and others in the background may also contribute largely. Many new operations have been made possible by advances in anesthesia, and these in turn depend upon engineers who have devised machines and chemists who have produced new drugs. Other operations are made possible by new materials, such as the alloys and plastics that are used to make artificial hip and knee joints.

Whenever practicable, new operations are tried on animals before they are tried on patients. This practice is particularly relevant to organ transplants. Surgeons themselves—not experimental physiologists—transplanted kidneys, livers, and hearts in animals before attempting these procedures on patients. Experiments on animals are of limited value, however, because animals do not suffer from all of the same maladies as do humans.

Many other developments in modern surgical treatment rest on a firm basis of experimentation, often first in animals but also in humans; among them are renal dialysis (the artificial kidney), arterial bypass operations, embryo implantation, and exchange transfusions. These treatments are but a few of the more dramatic of a large range of therapeutic measures that have not only provided patients with new therapies but also have led to the acquisition of new knowledge of how the body works. Among the research projects of the late 20th century is that of gene transplantation, which has the potential of providing cures for cancer and other diseases.

SCREENING PROCEDURES

Developments in modern medical science have made it possible to detect morbid conditions before a person actu-

The controlled clinical trial

New drug development

Hospital ethical committees

Animal experimentation

ally feels the effects of the condition. Examples are many: they include certain forms of cancer; high blood pressure; heart and lung disease; various familial and congenital conditions; disorders of metabolism, like diabetes; and acquired immune deficiency syndrome (AIDS). The consideration to be made in screening is whether or not such potential patients should be identified by periodic examinations. To do so is to imply that the subjects should be made aware of their condition and, second, that there are effective measures that can be taken to prevent their condition, if they test positive, from worsening. Such so-called specific screening procedures are costly since they involve large numbers of people. Screening may lead to a change in the life-style of many persons, but not all such moves have been shown in the long run to be fully effective. Although screening clinics may not be run by doctors, they are a factor of increasing importance in the preventive health service.

Periodic general medical examination of various sections of the population, business executives for example, is another way of identifying risk factors that, if not corrected, can lead to the development of overt disease.

(J.W.T./Ha.Sc.)

Medical education

Goals of medical education

Medical education is directed toward imparting to persons seeking to become physicians the knowledge and skills required for the prevention and treatment of disease and also for developing the methods and objectives appropriate to the study of the still unknown factors that produce disease or favour well-being. Among the goals of medical education is the production of physicians sensitive to the health needs of their country, capable of ministering to those needs, and aware of the necessity of continuing their own education. It therefore follows that the plan of education, the medical curriculum, should not be the same in all countries. Although there may be basic elements common to all, the details should vary from place to place and from time to time. Whatever form the curriculum takes, ideally it will be flexible enough to allow modification as circumstances alter, medical knowledge grows, and needs change.

Attention in this article is focused primarily on general medical education.

HISTORY OF MEDICAL EDUCATION

Although it is difficult to identify the origin of medical education, authorities usually consider that it began with the ancient Greeks' method of rational inquiry, which introduced the practice of observation and reasoning regarding disease. Rational interpretation and discussion, it is theorized, led to teaching and thus to the formation of schools such as that at Cos, where the Greek physician Hippocrates is said to have taught in the 5th century BC and originated the oath that became a credo for practitioners through the ages.

Later, the Christian religion greatly contributed to both the learning and the teaching of medicine in the West because it favoured not only the protection and care of the sick but also the establishment of institutions where collections of sick people encouraged observation, analysis, and discussion among physicians by furnishing opportunities for comparison. Apprenticeship training in monastic infirmaries and hospitals dominated medical education during the early Middle Ages. A medical school in anything like its present form, however, did not evolve until the establishment of the one at Salerno in southern Italy between the 9th and 11th centuries. Even there teaching was by the apprentice system, but an attempt was made at systemization of the knowledge of the time, a series of health precepts was drawn up, and a form of registration to practice was approved by the Holy Roman emperor Frederick II. During the same period, medicine and medical education were flourishing in the Muslim world at such centres as Baghdad, Cairo, and Córdoba.

With the rise of the universities in Italy and later in Cracow, Prague, Paris, Oxford, and elsewhere in western Europe, the teachers of medicine were in some measure

drawn away from the life of the hospitals and were offered the attractions and prestige of university professorships and lectureships. As a result, the study of medicine led more often to a familiarity with theories about disease than with actual sick persons. However, the establishment in 1518 of the Royal College of Physicians of London, which came about largely through the energies of Thomas Linacre, produced a system that called for examination of medical practitioners. The discovery of the circulation of the blood by William Harvey provided a stimulus to the scientific study of the processes of the body, bringing some deemphasis to the tradition of theory and doctrine.

Gradually, in the 17th and 18th centuries, the value of hospital experience and the training of the students' sight, hearing, and touch in studying disease were reasserted. In Europe, medical education began slowly to assume its modern character in the application of an increasing knowledge of natural science to the actual care of patients. There was also encouragement of the systematic study of anatomy, botany, and chemistry, sciences at that time considered to be the basis of medicine. The return to the bedside aided the hospitals in their long evolution from dwelling places of the poor, the diseased, and the infirm, maintained by charity and staffed usually by religious orders, into relatively well-equipped, well-staffed, efficient establishments that became available to the entire community and were maintained by private or public expense.

It was not until the mid-19th century, however, that an ordered pattern of science-oriented teaching was established. This pattern, the traditional medical curriculum, was generally adopted by Western medical schools. It was based upon teaching, where the student mostly listens, rather than learning, where the student is more investigative. The clinical component, largely confined to hospitals (charitable institutions staffed by unpaid consultants), was not well organized. The new direction in medical education was aided in Britain by the passage of the Medical Act of 1858, which has been termed the most important event in British medicine. It established the General Medical Council, which thenceforth controlled admission to the medical register and thus had great powers over medical education and examinations. Further interest in medicine grew from these advances, which opened the way for the discoveries of Louis Pasteur, which showed the relation of microorganisms to certain diseases, Joseph Lister's application of Pasteur's concepts to surgery, and the studies of Rudolf Virchow and Robert Koch in cellular pathology and bacteriology.

In the United States, medical education was greatly influenced by the example set in 1893 by the Johns Hopkins Medical School in Baltimore. It admitted only college graduates with a year's training in the natural sciences. Its clinical work was superior because the school was supplemented by the Johns Hopkins Hospital, created expressly for teaching and research carried on by members of the medical faculty. The adequacy of medical schools in the United States was improved after the Carnegie Foundation for the Advancement of Teaching published in 1910 a report by the educator Abraham Flexner. In the report, which had an immediate impact, he pointed out that medical education actually is a form of education rather than a mysterious process of professional initiation or apprenticeship. As such, it needs an academic staff, working full-time in their departments, whose whole responsibility is to their professed subject and to the students studying it. Medical education, the report further stated, needs laboratories, libraries, teaching rooms, and ready access to a large hospital, the administration of which should reflect the presence and influence of the academic staff. Thus the nature of the teaching hospital was also influenced. Aided by the General Education Board, the Rockefeller Foundation, and a large number of private donors, U.S. and Canadian medical education was characterized by substantial improvements from 1913 to 1929 in such matters as were stressed in the Flexner report.

MODERN PATTERNS OF MEDICAL EDUCATION

As medical education developed after the Flexner report was published, the distinctive feature was the thorough-

The Royal College of Physicians

Advances at Johns Hopkins Medical School

ness with which theoretical and scientific knowledge were fused with what experience teaches in the practical responsibility of taking care of human beings. Medical education eventually developed into a process that involved four generally recognized stages: premedical, undergraduate, postgraduate, and continuing education.

Premedical education and admission to medical school. In the United States, Britain, and the Commonwealth countries, generally, medical schools are inclined to limit the number of students admitted so as to increase the opportunities for each student. In western Europe, South America, and most other countries, no exact limitation of numbers of students is in effect, though there is a trend toward such limitation in some of the western European schools. Some medical schools in North America have developed ratios of teaching staff to students as high as 1 to 1 or 1 to 2, in contrast with 1 teacher to 20 or even 100 students in certain universities in other countries. The number of students applying to medical school greatly exceeds the number finally selected in most countries.

Requirements to enter medical school, of course, vary from country to country, and in some countries, such as the United States, from university to university. Generally speaking, in Western universities, there is a requirement for a specified number of years of undergraduate work and passing of a test, possibly state regulated, and a transcript of grades. In the United States entry into medical school is highly competitive, especially in the more prestigious universities. Stanford University, for instance, accepts only about 5 percent of its applicants. Most U.S. schools require the applicant to take the Medical College Admission Test, which measures aptitude in medically related subjects. Other requirements may include letters of recommendation and a personal interview. Many U.S. institutions require a bachelor's degree or its equivalent from an undergraduate school. A specific minimum grade point average is not required, but most students entering medical school have between an A and a B average.

The premedical courses required in most countries emphasize physics, chemistry, and biology. These are required in order to make it possible to present subsequently courses in anatomy, physiology, biochemistry, and pharmacology with precision and economy of time to students prepared in scientific method and content. Each of the required courses includes laboratory periods throughout the full academic year. Student familiarity with the use of instruments and laboratory procedures tends to vary widely from country to country, however.

Undergraduate education. The medical curriculum also varies from country to country. Most U.S. curriculums cover four years; in Britain five years is normal. The early part of the medical school program is sometimes called the preclinical phase. Medical schools usually begin their work with the study of the structure of the body and its formation: anatomy, histology, and embryology. Concurrently, or soon thereafter, come studies related to function—*i.e.*, physiology, biochemistry, pharmacology, and, in many schools, biophysics. After the microscopic study of normal tissues (histology) has begun, the student is usually introduced to pathological anatomy, bacteriology, immunology, parasitology—in short, to the agents of disease and the changes that they cause in the structure and function of the tissues. Courses in medical psychology, biostatistics, public health, alcoholism, biomedical engineering, emergency medicine, ethical problems, and other less traditional courses are becoming more common in the first years of the medical curriculum.

The two or more clinical years of an effective curriculum are characterized by active student participation in small group conferences and discussions, a decrease in the number of formal lectures, and an increase in the amount of contact with patients in teaching hospitals and clinics.

Clinical work begins with general medicine and surgery and goes on to include the major clinical specialties, including obstetrics and gynecology, pediatrics, disorders of the eye, ear, nose, throat, and skin, and psychiatry. The student works in the hospital's outpatient, emergency, and radiology departments, diagnostic laboratories, and surgical theatres. The student also studies sciences

closely related to medicine, such as pathology, microbiology, hematology, immunology, and clinical chemistry and becomes familiar with epidemiology and the methods of community medicine. Some knowledge of forensic (legal) medicine is also expected. During the clinical curriculum many students have an opportunity to pursue a particular interest of their own or to enlarge their clinical experience by working in a different environment, perhaps even in a foreign country—the so-called elective period. Most students find clinical work demanding, usually requiring long hours of continuous duty and personal commitment.

In the United States after satisfactory completion of a course of study in an accredited medical school the degree of doctor of medicine (M.D.) or doctor of osteopathy (D.O.) is conferred. In Britain and some of the other Commonwealth countries the academic degree conferred after undergraduate studies are completed is bachelor of medicine and of surgery (or *chirurgery*), M.B., B.S. or M.B., ChB. Only after further study is the M.D. degree given. Similar degrees are conferred in other countries, although they are not always of the same status.

Postgraduate education. On completion of medical school, the physician usually seeks graduate training and experience in a hospital under the supervision of competent clinicians and other teachers. In Britain a year of resident hospital work is required after qualification and before admission to the medical register. In North America, the first year of such training has been known as an internship, but it is no longer distinguished in most hospitals from the total postgraduate period, called residency. After the first year physicians usually seek further graduate education and training to qualify themselves as specialists or to fulfill requirements for a higher academic degree. Physicians seeking special postgraduate degrees are sometimes called fellows.

Continuing education. The process by which physicians keep themselves up-to-date is called continuing education. It consists of courses and training opportunities of from a few days to several months in duration, designed to enable physicians to learn of new developments within their special areas of concern. Physicians also attend medical and scientific meetings, national and international conferences, discussion groups, and clinical meetings, and they read medical journals and other materials, all of which serve to keep them aware of progress in their chosen field. Although continuing education is not a formal process, organizations designed to promote continuing education have become common. In the United States the Accreditation Council for Continuing Medical Education was formed in 1985, and some certifying boards of medical specialties have stringent requirements for continuing education.

The quality of medical education is supervised in many countries by councils appointed by the profession as a whole. In the United States these include the Council on Medical Education and the Liaison Committee on Medical Education, both affiliates of the American Medical Association, and the American Osteopathic Association. In Britain the statutory body is the General Medical Council, most of whose members are from the profession, although only a minority of the members are appointed by it. In other countries medical education may be regulated by an office or ministry of public instruction with, in some cases, the help of special professional councils.

Medical school faculty. As applied to clinical teachers the term full-time originally implied an educational ideal: that a clinician's salary from a university should be large enough to relieve him of any reason for seeing private patients for the sake of supplementing his salary by professional fees. Full-time came to be applied, however, to a variety of modifications; it could mean that a clinical professor might supplement his salary as a teacher up to a defined maximum, might see private patients only at his hospital office, or might see such patients only a certain number of hours per week. The intent of full-time has always been to place the teacher's capacities and strength entirely at the service of his students and the patients entrusted to his care as a teacher and investigator.

Courses in the medical sciences have commonly followed the formula of three hours of lectures and six to nine

Limitation
of number
of students

Internship
and
residency

The
clinical
years

hours of laboratory work per week for a three-, six-, or nine-month course. Instruction in clinical subjects, though retaining the formal lecture, have tended to diminish the time and emphasis allowed to lectures in favour of experience with and attendance on patients. Nonetheless, the level of lecturing and formal presentation remains high in some countries.

REQUIREMENTS FOR PRACTICE

Graduation from medical school and postgraduate work does not always allow the physician to practice. In the United States, licensure to practice medicine is controlled by boards of licensure in each state. The boards set and conduct examinations of applicants to practice within the state, and they examine the credentials of applicants who want licenses earned in other states to be accepted in lieu of examination. The National Board of Medical Examiners holds examinations leading to a degree that is acceptable to most state boards. National laws regulating professional practice cannot be enacted in the United States. In Canada the Medical Council of Canada conducts examinations and enrolls successful candidates on the Canadian medical register, which the provincial governments accept as the main requirement for licensure. In Britain the medical register is kept by the General Medical Council, which supervises the licensing bodies; unregistered practice, however, is not illegal. In some European countries graduation from a state-controlled university or medical school in effect serves as a license to practice; the same is true for Japan. (For further details of licensure, see below *Legal aspects of medicine*.)

ECONOMIC ASPECTS

The income of a medical school is derived from four principal sources: (1) tuition and fees, (2) endowment income or appropriation from the government (taxation), (3) gifts from private sources, and (4) donation of teachers' services. Tuition or student fees are large in most English-speaking countries (except in U.S. state universities) and relatively small throughout the rest of the world. Tuition in most American schools, however, rarely makes up more than a small part of total operating expenses. The total cost of maintaining a medical school, if prorated among the students, would produce a figure many times greater than the tuition or other charges paid by each student. The costs of operating medical schools in the United States increased by about 30 times between the late 1950s and the mid-1980s.

The expenses of medical education fall into two groups: those of the instruction given in the medical sciences and those connected with hospital teaching. In the medical sciences the costs of building maintenance, laboratory equipment and supplies, research expenses, salaries of teachers, and wages of employees are heavy but comparable to those in other departments of a university. In the clinical sub-

jects all expenses in connection with the care of patients usually are considered as hospital expenses and are not carried on the medical school budget, which is normally reserved for the expenses of teaching and research. Here the heavy expenses are salaries of clinical teachers and the cost of studying cases of illness with a thoroughness appropriate to their use as teaching material.

To a considerable degree in free-market countries, the cost of securing an adequate medical education has tended to exclude the student whose family cannot contribute a large share of tuition and living expenses for four to 10 years. This difficulty is offset in some medical schools by loan funds and scholarships, but these aids are commonly offered only in the second or subsequent years. In Britain scholarships and maintenance grants are available through state and local educational authority funds, so that an individual can secure a medical education even though the parents may not be able to afford its cost.

SCIENTIFIC AND INTERNATIONAL ASPECTS

Medical education has the double task of passing on to students what is known and of attacking what is still unknown. The cost of medical research is borne by only a few; the benefits are shared by many. There are countries whose citizens are too poor to support physicians or to use them, countries that can support a few physicians but are too poor to maintain a good medical school, countries that can maintain medical schools where what is known can be taught but where no research can be carried out, and a few countries in which teaching and research in medicine can be carried on to the great advantage of the world at large.

A medical school having close geographical as well as administrative relationships with the rest of the university of which it forms a part usually profits by this intimate and easy contact. Medicine cannot wisely be separated from the biological sciences, and it continues to gain immensely from chemistry, physics, mathematics, and psychology, as well as from modern technology. The social sciences contribute by making physicians aware of the need for better distribution of medical care. Contact with teachers and the advancing knowledge in other faculties also may have a corollary effect in advancing medicine.

With the development of the World Health Organization (WHO) and the World Medical Association after World War II, there has been increasing international interest in medical education. WHO conducts a regular program for aiding countries in the development and expansion of their educational facilities. World War II showed the advantages and economy derived from satisfactory systems of medical education: defects and diseases were more widely and accurately detected among recruits than ever before, health and morale were effectively maintained among combatants, and disease and battle injuries were effectively treated. (Al.Gr./E.L.T./Ha.Sc.)

WHO

MAJOR MEDICAL INSTITUTIONS

Hospitals

A hospital is an institution that is built, staffed, and equipped for the identification (diagnosis) of disease: for the treatment, both medical and surgical, of the sick and the injured; and for their housing during this process. The modern hospital also often serves as a centre for investigation and for teaching. To better serve the wide-ranging needs of the community, the modern hospital has often developed outpatient facilities, as well as emergency, psychiatric, and rehabilitation services.

Hospitals have long existed in every civilized country. The developing countries, which contain a large proportion of the world's population, do not have enough hospitals, equipment, and trained staff, and, by the standards of the industrialized countries, the hospitals that do exist are poorly equipped to handle the volume of persons who need care. These persons, then, do not always receive the benefits of modern medicine, public health

measures, or hospital care, and they generally have lower life expectancies.

In the developed countries the hospital as an institution is becoming more complex as modern technology increases the range of diagnostic capabilities and expands the possibilities for treatment. As a result of the greater range of services and the more involved treatment and surgery available, the ratio of staff to patient has increased and a more highly trained staff is required. During recent years a combination of medicine and engineering has produced a vast array of new instrumentation, much of which requires a hospital setting for its use. Hospitals thus are becoming more expensive to run, and health service administrators are increasingly concerned with the question of cost-effectiveness.

HISTORY OF HOSPITALS

As early as 4000 bc religions identified certain of their deities with healing. The temples of Saturn, and later of

Licensure to practice

Sources of medical school income

Hospitals
in
antiquity

Asclepius in Asia Minor, were recognized as healing centres. Brahmanic hospitals were established in Sri Lanka as early as 431 BC, and King Aśoka established a chain of hospitals in Hindustān about 230 BC. Around 100 BC the Romans established hospitals (*valetudinaria*) for the treatment of their sick and injured soldiers; their care was important because it was upon the integrity of the legions that the power of Rome was based.

It can be said, however, that the modern concept of a hospital dates from AD 331 when Constantine, having been converted to Christianity, abolished all pagan hospitals and thus created the opportunity for a new start. Until that time disease had isolated the sufferer from the community. The Christian tradition emphasized the close relationship of the sufferer to his fellow man, upon whom rested the obligation for care. Illness thus became a matter for the Christian church.

Around AD 370 St. Basil of Caesarea established a religious foundation in Cappadocia that included a hospital, an isolation unit for those suffering from leprosy, and buildings to house the poor, the elderly, and the sick. Following this example similar hospitals were later built in the eastern part of the Roman Empire. Another notable foundation was that of St. Benedict at Monte Cassino, founded early in the 6th century, where the care of the sick was placed above and before every other Christian duty. It was from this beginning that one of the first medical schools in Europe ultimately grew at Salerno and was of high repute by the 11th century. This example led to the establishment of similar monastic infirmaries in the western part of the empire.

The Hôtel-Dieu of Lyon was opened in 542 and the Hôtel-Dieu of Paris in 660. In these hospitals more attention was given to the well-being of the patient's soul than to curing bodily ailments. The manner in which monks cared for their own sick became a model for the laity. The monasteries had an *infirmatorium*, a place to which their sick were taken for treatment. The monasteries had a pharmacy and frequently a garden with medicinal plants. In addition to caring for sick monks, the monasteries opened their doors to pilgrims and to other travelers.

Religion continued to be the dominant influence in the establishment of hospitals during the Middle Ages. The growth of hospitals accelerated during the Crusades, which began at the end of the 11th century. Pestilence and disease were more potent enemies than the Saracens in defeating the crusaders. Military hospitals came into being along the traveled routes; the Knights Hospitalers of the Order of St. John in 1099 established in the Holy Land a hospital that could care for some 2,000 patients. It is said to have been especially concerned with eye disease, and may have been the first of the specialized hospitals. This order has survived through the centuries as the St. John's Ambulance Corps.

Throughout the Middle Ages, but notably in the 12th century, the number of hospitals grew rapidly in Europe. The Arabs established hospitals in Baghdad and Damascus and in Córdoba in Spain. Arab hospitals were notable for the fact that they admitted patients regardless of religious belief, race, or social order. The Hospital of the Holy Ghost, founded in 1145 at Montpellier in France, established a high reputation and later became one of the most important centres in Europe for the training of doctors. By far the greater number of hospitals established during the Middle Ages, however, were monastic institutions under the Benedictines, who are credited with having founded more than 2,000.

The Middle Ages also saw the beginnings of support for hospital-like institutions by secular authorities. Toward the end of the 15th century many cities and towns supported some kind of institutional health care: it has been said that in England there were no less than 200 such establishments that met a growing social need. This gradual transfer of responsibility for institutional health care from the church to civil authorities continued in Europe after the dissolution of the monasteries in 1540 by Henry VIII, which put an end to hospital building in England for some 200 years.

The loss of monastic hospitals in England caused the

secular authorities to provide for the sick, the injured, and the handicapped, thus laying the foundation for the voluntary hospital movement. The first voluntary hospital in England was probably established in 1718 by Huguenots from France and was closely followed by the foundation of such London hospitals as the Westminster Hospital in 1719, Guy's Hospital in 1724, and the London Hospital in 1740. Between 1736 and 1787 hospitals were established outside London in at least 18 cities. The initiative spread to Scotland where the first voluntary hospital, the Little Hospital, was opened in Edinburgh in 1729.

The first hospital in North America was built in Mexico City in 1524 by Cortés; the structure still stands. The French established a hospital in Canada in 1639 at Quebec city, the Hôtel-Dieu du Précieux Sang, which is still in operation although not at its original location. In 1644 Jeanne Mance, a French noblewoman, built a hospital of ax-hewn logs on the island of Montreal; this was the beginning of the Hôtel-Dieu de St. Joseph, out of which grew the order of the Sisters of St. Joseph, now considered to be the oldest nursing group organized in North America. The first hospital in the territory of the present-day United States is said to have been a hospital for soldiers on Manhattan Island, established in 1663.

The early hospitals were primarily almshouses, one of the first of which was established by William Penn in Philadelphia in 1713. The first incorporated hospital in America was the Pennsylvania Hospital, in Philadelphia, which obtained a charter from the crown in 1751.

THE MODERN HOSPITAL

Hospitals may be compared and classified in various ways: by ownership and control, by type of service rendered, by length of stay, by size, or by facilities and organization provided. Terms in general use include the general as distinct from the special hospital, the short-stay hospital, and the long-term hospital.

Comparison. Hospitals may be compared by the number of beds they contain. Modern hospitals tend to be small, rarely exceeding 800 beds, which is thought to be the largest number that can be governed satisfactorily from a single administrative unit, yet not too large to retain a corporate unity.

Another index is the average bed-occupancy, that is, the percentage of available beds actually occupied per day or per month. In Europe bed-occupancy may be higher in the cold winter months, which bring more respiratory disease. In developing countries the bed-occupancy is often more than 100 percent: that is to say there are more patients in the hospital than there are beds for them.

The amount of time that a patient spends in a hospital bed, the bed-stay, is another important index and depends on the nature of the hospital. In an acute-care hospital the bed-occupancy will be low. In hospitals catering to the more chronically ill, the average bed-stay probably will be higher. There are even significant variations between units in the same hospital doing the same kind of work. In hospitals in developing countries, the average bed-stay is much shorter than in Europe.

Ownership, control, and financing. *Ownership and control.* In most countries outside of North America nearly all hospitals are owned and operated by the government. In Great Britain, except for a small number run by religious orders or serving special groups, most hospitals are within the National Health Service. The local hospital management committee answers directly to the regional hospital board and ultimately to the Department of Health and Social Security. In the United States and Canada most hospitals are nonprofit and are neither owned nor operated by governmental agencies. Many of them are associated with universities; others were founded by religious groups or by public-spirited individuals. Mental hospitals traditionally have been the responsibility of the state governments, while military and veterans hospitals have been provided by the federal government. In addition, there are a number of municipal and county general hospitals.

Financing. Almost universally, hospital construction costs are met at least in some part by governmental contributions. Operating costs are taken care of in a variety

Early
hospitals
in North
America

of ways. Ultimately, a substantial portion of the expenses not covered by private endowments or gifts is met by contributions from the general funds of some unit of government or out of funds collected by insurance carriers from subscribers. In countries in which hospital insurance is not universal or its coverage is incomplete, some of the operating costs are met by charges on uninsured or inadequately insured patients.

Carriers
of hospital
insurance

The carriers of the hospital insurance in a particular country may be governmental agencies, private corporations or agencies, or both. In Britain, for example, under the National Insurance Act, the government is the carrier. All persons who have reached the minimum school-leaving age and are not full-time students, beyond the age of retirement, in prison, or receivers of benefits from the insurance and who do not have less than a certain minimum income are contributors under the plan whether employed by others, self-employed, or nonemployed. Employers also contribute.

In the United States persons who are employed by others or are self-employed make compulsory contributions toward a form of national hospital insurance, Medicare, which pays a large portion of the hospital costs of persons at the age of 65 or over. Employers then make matching payments. A majority of the persons ineligible, by reason of age, for benefits under the Medicare program are enrolled in some other form of hospital insurance, such as one of the plans offered by the commercial insurance companies or one of the independent plans, including community and community-controlled plans and those operated by unions, employers, welfare funds, and private medical clinics.

In the United States, even with federal participation under Medicare and Medicaid (a program for persons under 65 who are unable to pay), the payment for health care services on an insurance basis, either voluntary or governmental, is considerably less advanced than it is in many other parts of the world. In Europe, particularly, the financial support of services in hospitals tends to be much more collectivized. Less than 10 percent of the costs of hospital operation in Europe is covered by payments made directly by patients. Details vary somewhat from country to country: in the United Kingdom, for example, the funds for total hospital operation are appropriated by the Ministry of Social Security to each regional hospital board, which in turn distributes them to the local hospital groups. In Sweden, however, approximately 90 percent of hospital operating costs are provided by local or provincial units of government from public revenue; the remaining 10 percent of the costs comes from payments made by insurance funds on behalf of the patient. In general, in the majority of European countries, hospital operating costs are paid out of insurance funds; such is the case in France, Italy, The Netherlands, Norway, and elsewhere.

The general hospital. General hospitals are general in the sense that they admit all types of medical and surgical cases, and they concentrate on patients with acute illness needing relatively short-term care. A community general hospital with about 200 beds has an organized medical staff, a professional nursing staff, and diagnostic equipment. In addition to the essential services relating to patient care, it has a pharmacy, a laboratory, X-ray and physical therapy departments, possibly a maternity division (ordinarily including a nursery and a delivery room), operating rooms, recovery rooms, an outpatient department, and an emergency department.

The
community
general
hospital

In a somewhat larger hospital there may be additional facilities: dental services, a nursery for premature infants, an organ bank for use in transplantation, a department of renal dialysis (removal of wastes from the blood by passing it through semipermeable membranes, as in the artificial kidney), equipment for inhalation therapy, an intensive-care unit, a volunteer-services department, and, possibly, a home-care program. The complexity of the general hospital, then, reflects the advances made after World War II, including the use of antibiotics, a vast new array of laboratory procedures, new surgical techniques, new materials and equipment for radiation therapy, and an increased emphasis on physical therapy and rehabilitation.

The legally constituted governing body of the hospital, with full responsibility for the conduct and efficient management of the hospital, is usually a hospital board. The board establishes policy and, on the advice of a medical advisory board, appoints a medical staff and an administrator. It exercises control over expenditures and has the responsibility for maintaining professional standards.

The administrator is the chief executive officer of the hospital and is responsible to the board. In a large hospital there are many separate departments, each of which is controlled by a department head. The largest department in any hospital is nursing, followed by the dietary department and housekeeping. Other departments that are important to the functioning of the hospital include laundry, engineering, stores, purchasing, accounting, pharmacy, physical therapy, social service, pathology, X-ray, and medical records.

The medical staff is organized into such departments as surgery, medicine, obstetrics, and pediatrics. The degree of departmentalization of the medical staff depends on the specialization of its members and not primarily on the size of the hospital, although there is usually some correlation between the two. The chiefs of the medical-staff departments, along with the chiefs of radiology and pathology, make up the medical advisory board, which usually holds monthly meetings on medical-administrative matters. The professional work of the individual staff members is reviewed by medical-staff committees. In a large hospital the committees may report to the medical advisory board; in a smaller hospital, to the medical staff directly, at regular staff meetings.

Organiza-
tion of
medical
staff

General hospitals often also have a formal or informal role as teaching institutions. When formally designed as such, teaching hospitals are affiliated with undergraduate and postgraduate medical education at a university, and they provide up-to-date and often specialized therapeutic measures and facilities unavailable elsewhere in the region. As teaching hospitals have become more specialized, general hospitals have become more involved in providing general clinical training to medical students.

Specialized health- and medical-care facilities. Hospitals that specialize in one type of illness or one type of patient can be found in Europe and in North America, although, except in large university centres where postgraduate teaching is carried out on a large scale, the special hospital is increasingly becoming a department of the general hospital. Changing conditions or modes of treatment have lessened the need or reduced the number of some types of specialized institutions; this may be seen in the cases of tuberculosis, leprosy, and mental hospitals.

Tuberculosis and leprosy hospitals. Between 1880 and 1940 tuberculosis hospitals provided rest, relaxation, special diets, and fresh air, and even if the tuberculosis was in an early stage, a stay of more than two years was thought necessary to effect a healing of the disease; a permanent cure was not considered entirely feasible. Today the use of antibiotics, along with advances in chest surgery and routine X-ray programs, has meant that the treatment of tuberculosis need not be carried out in a specialized facility.

Leprosy has been known for centuries to be contagious. Lazar houses were established throughout Europe in the Middle Ages to isolate those with leprosy, at that time a common disease, from the community. In the 14th century there may have been some 7,000 leper houses in France alone, and some of the earliest hospitals in England were established for lepers. Although it is now rare in Europe, leprosy is still common in many parts of the world. The purpose of the modern leprosarium is not so much isolation as it is treatment. The chronic form of the disease is treated by surgical correction of deformities, occupational therapy, rehabilitation, and sheltered living in associated villages. Acute leprosy is treated in general hospitals, clinics, and dispensaries.

Mental hospitals. Psychiatric patients traditionally have been housed in long-stay mental hospitals, formerly called asylums, although the majority of large general hospitals now have a psychiatric unit. In the past few decades the hospital stay of many persons with chronic mental disease has been shortened by modern medication and better

Changes in treatment of mental patients

understanding on the part of the public. Mental patients often may participate in many activities, first within the hospital setting and later in the community, either with trial visits at home or with placement of selected patients in foster homes. Effort is now made with appropriate medication and the judicious use of the support services to get the patient home and in the care of the family. Even those mentally handicapped persons who require custodial care are no longer isolated from contact with their relatives and friends.

Long-stay institutions. Historically the long-stay institution was a place for the elderly, the infirm, and those with chronic irreversible and disabling disorders, especially if the patients were indigent. Medical and nursing care was minimal. Today the long-stay hospital has a more active role in health care. Many are well staffed and well equipped to help a patient prepare to live at home or with a member of the family. Long-stay hospitals represent a significant extension of the hospital health-care system, helping to conserve expensive facilities for the acutely ill and improving the prospects of the chronically disabled.

Private hospitals. Throughout Europe, as well as in North America, Australia, and New Zealand, there are small private hospitals, often called nursing homes in Britain, many of which until recently were able to provide little more than accommodation and simple nursing, the patient being under the care of a general practitioner or of a visiting consultant. Medical practice in the towns of developing countries is characterized by a proliferation of many small private hospitals, usually owned by doctors, that have developed to meet the widespread need for hospital care not otherwise available.

Another method of providing health care in a hospital for those able to pay for it, both in industrial and developing countries, is the provision of a limited number of beds for private patients within a large general hospital otherwise financed to some degree by public funds. In the United Kingdom and, for example, in West Africa, these so-called amenity beds usually form part of the ward unit, the patient being required to pay for certain amenities such as a measure of privacy, unrestricted visiting, attractively served food, and a more liberal routine. Alternatively, many large general hospitals are able to offer much more costly accommodations in so-called private blocks—that is, in a part of the hospital specially designed and equipped for private patients. Patients in a private block pay a large portion of the total cost of their medical care, including that of surgery.

More recent is the development of the wholly independent private hospital run by a company or business consortium. Many of these privately funded hospitals are able to provide most or all of the services of a general hospital, including constant medical care and a first-rate nursing service, although such facilities are costly.

The hospice. Historically a hospice was a guest house intended for pilgrims and was often closely connected with a monastery and supervised by monks. From the beginning it had a strong religious connection and exemplified the Christian insistence on compassion and care for the aged, the infirm, the needy, and the ill. In modern Britain the hospice movement developed gradually from its beginning in 1905, when the Sisters of Charity founded the St. James Hospice in London. The St. Christopher Hospice, also in London, founded in 1967, soon became known for its peaceful environment and expert medical and nursing care. The hospice movement has spread throughout Britain, North America, Australia, and Europe.

Mission hospitals. The spread of Western medicine and the founding of hospitals in the developing countries can be attributed in large part to the influence of the medical missionary. The establishment of mission hospitals gained momentum gradually in the second half of the 19th century. By the second half of the 20th century, however, this steady growth had already dwindled, since all but a few of the hospitals and dispensaries founded during that hundred years had been absorbed into the native health-care system. The Christian missionaries had a great influence on the creation of centres of Western medicine in many developing countries and in promulgating the concept of

a hospital in which health care would be centralized and organized for the benefit of the ill and injured, many of whom would not otherwise have survived. The medical missionaries also promoted the idea and the ideals of nursing as a profession for native men and women.

Apart from its religious associations, a mission hospital functions as a general hospital in the sense that it admits all who need hospital care. A number of mission hospitals, however, have been devoted to specific diseases—for example, leprosy and diseases of the eyes. Perhaps the most important contribution made by mission hospitals is in the enormous numbers of persons, particularly women and children, who have been treated as outpatients.

Extended health care. With the advance in medical science and the ever-increasing cost of hospital operations, the progressive-care concept is more attractive, both for outpatient and inpatient care. Progressive care can be divided into five categories: (1) intensive care, (2) intermediate care, (3) self-care, (4) long-term care, and (5) organized home-care programs. Two of the categories, self-care and home-care programs, are relatively new departures from past practice and deserve special attention.

Self-care facilities are organized into a separate unit in which ambulatory patients who require only diagnostic or convalescent care are given accommodations similar to those of a hotel. The patients are free to wear street clothes and to go to the hospital cafeteria. Such a ward or wing of a general hospital requires much less costly equipment than the intensive- or intermediate-care units and can be staffed with far fewer nurses and aides.

Home-care programs are for patients who need some care but not all of the treatment facilities of a hospital. The patients are provided with a range of individualized medical, nursing, social, and rehabilitative services in their own homes, coordinated through one central agency. Patients can be considered ready for home care when the following criteria are met: (1) diagnosis and a plan for treatment have been established; (2) inpatient hospital facilities are no longer required for proper care; (3) no more than two visits per week by physicians are required; (4) the nursing service has found that the physical environment of the home is such that the patient receives adequate care; (5) the patient is too ill to visit an outpatient clinic but does not need hospital care; (6) the family environment would have a therapeutic effect, and family members or others can be taught to provide the necessary care; and (7) the family and the patient prefer that care be provided at home. Even though home care conserves expensive acute-care beds, and although most patients on home care do as well as or better than expected, home-care programs have not been widely adopted.

Regional planning. Sweden and the former Soviet Union provide examples of advanced planning in the integration of hospital networks into coordinated health services. In both nations the government was charged with the responsibility of providing health care to all citizens. In Sweden financing is in part by compulsory health insurance.

Sweden is divided into health service regions; each region includes several counties and has a central hospital. Each county within a region has a county hospital with up to 1,000 beds and with specialized and outpatient facilities to serve a population of about 300,000. The counties, in turn, are divided into districts, each of which has a population of about 75,000 and is served by a district hospital, which usually has 300 or more beds. Smaller communities have health centres or ambulatory service centres that are not administered as part of the hospital system.

The Soviet Union took a somewhat different approach. In its thinly populated rural areas, general hospitals, called *uchastok* hospitals, served populations as small as 2,000 to 15,000 persons. These 15- to 100-bed general hospitals occupied the same premises and employed the same staff as general clinics (polyclinics) that provided general and specialized care. The hospital-clinic staff included a general physician, a surgeon, and a dental surgeon. Some larger *uchastok* centres also had a radiologist and a pathologist.

The next larger hospitals, the district hospitals, had 250–500 beds and usually had divisions for surgical, medical, obstetric, and pediatric services and provided care for in-

Categories of progressive care

Private patients

fectious diseases; some also included departments for eye, ear, nose, and throat disorders and for orthopedic surgery. Patients who could not be treated adequately in the district hospitals were referred to the next higher level, the regional hospital, which served a population of 1,000,000–5,000,000 people and contained up to 1,250 beds.

The republic hospital occupied the highest level in the Soviet system. Such a hospital, or complex of hospitals, served as a referral centre and had the responsibility of undergraduate medical education. Some were also associated with one or more research institutes.

Regional planning in North America is less advanced. One regional pattern is a satellite system, centred on a metropolis and applying the principle of progressive patient care. The system is focused on the efficient provision of comprehensive health care to the residents of the region. Less serious cases are handled in the outer, more accessible health facilities of the system; the more serious are referred to the inner hospitals of the ring or to the research and teaching hospital at the core.

The term metropolitan planning council is often used to denote an advisory planning group that coordinates services among member hospitals in a metropolitan area and decides such questions as where more beds are to be added. In North America, however, most hospitals are not government-operated, and it is difficult to achieve close cooperation among voluntary groups.

(W.D.P./Ha.Sc./Ed.)

Public health services

Public health has been defined as the art and science of preventing disease, prolonging life, and promoting physical and mental health, sanitation, personal hygiene, control of infection, and organization of health services. From the normal human interactions involved in dealing with the many problems of social life, there has emerged a recognition of the importance of community action in the promotion of health and the prevention and treatment of disease; this is expressed in the concept of public health.

Comparable terms for public health medicine are social medicine and community medicine; the latter has been widely adopted in the United Kingdom, and the practitioners are called community physicians. The practice of public health draws heavily on medical science and philosophy and concentrates especially on manipulating and controlling the environment for the benefit of the public. It is concerned therefore with housing, water supplies, and food. Noxious agents can be introduced into these through farming, fertilizers, inadequate sewage disposal and drainage, construction, defective heating and ventilating systems, machinery, and toxic chemicals. Public health medicine is part of the greater enterprise of preserving and improving the public health. Community physicians cooperate with such diverse groups as architects, builders, sanitary and heating and ventilating engineers, factory and food inspectors, psychologists and sociologists, chemists, physicists, and toxicologists. Occupational medicine is concerned with the health, safety, and welfare of persons in the workplace. It may be viewed as a specialized part of public health medicine since its aim is to reduce the risks in the environment in which persons work.

The venture of preserving, maintaining, and actively promoting public health requires special methods of information-gathering (epidemiology) and corporate arrangements to act upon significant findings and put them into practice. Statistics collected by epidemiologists attempt to describe and explain the occurrence of disease in a population by correlating factors such as diet, environment, radiation, or cigarette smoking with the incidence and prevalence of disease. The government, through laws and regulations, creates agencies to oversee and formally inspect such things as water supplies, food processing, sewage treatment, drains, air contamination, and pollution. Governments also are concerned with the control of epidemic infections by means of enforced quarantine and isolation—for example, the health control that takes place at seaports and airports in an attempt to assure that infectious diseases are not brought into a country.

This section traces the historical development of public health, beginning in ancient times and emphasizing how various public health concepts have evolved. It outlines the organizational and administrative methods of handling these problems in the developed and the developing countries of the world. Special attention is given to the developing countries and to how the health problems, limitations of resources, education of health personnel, and other factors must be taken into account in designing health service systems. Finally, there are descriptions of the most recent developments in public health, together with some indications of the problems still to be solved.

HISTORY OF PUBLIC HEALTH

Beginnings in antiquity. Most of the world's primitive people have practiced cleanliness and personal hygiene, often for religious reasons, including, apparently, a wish to be pure in the eyes of their gods. The Old Testament, for example, has many adjurations and prohibitions about clean and unclean living. Religion, law, and custom were inextricably interwoven. For thousands of years primitive societies looked upon epidemics as divine judgments on the wickedness of mankind. The idea that pestilence is due to natural causes, such as climate and physical environment, however, gradually developed. This great advance in thought took place in Greece during the 5th and 4th centuries BC and represented the first attempt at a rational, scientific theory of disease causation. The association between malaria and swamps, for example, was established very early (503–403 BC), even though the reasons for the association were obscure. In the book *Airs, Waters, and Places*, thought to have been written by Hippocrates in the 5th or 4th century BC, the first systematic attempt was made to set forth a causal relationship between human diseases and the environment. Until the new sciences of bacteriology and immunology emerged well into the 19th century, this book provided a theoretical basis for the comprehension of endemic disease (that persisting in a particular locality) and epidemic disease (that affecting a number of people within a relatively short period).

The Middle Ages. In terms of disease, the Middle Ages can be regarded as beginning with the plague of 542 and ending with the Black Death (bubonic plague) of 1348. Diseases in epidemic proportions included leprosy, bubonic plague, smallpox, tuberculosis, scabies, erysipelas, anthrax, trachoma, sweating sickness, and dancing mania (see INFECTIOUS DISEASES). The isolation of persons with communicable diseases first arose in response to the spread of leprosy. This disease became a serious problem in the Middle Ages and particularly in the 13th and 14th centuries.

The Black Death reached the shores of southern Europe from the Middle East in 1348 and in three years swept throughout Europe. The chief method of combating plague was to isolate known or suspected cases as well as persons who had been in contact with them. The period of isolation at first was about 14 days and gradually was increased to 40 days. Stirred by the Black Death, public officials created a system of sanitary control to combat contagious diseases, using observation stations, isolation hospitals, and disinfection procedures. Major efforts to improve sanitation included the development of pure water supplies, garbage and sewage disposal, and food inspection. These efforts were especially important in the cities, where people lived in crowded conditions in a rural manner with many animals around their homes.

During the Middle Ages a number of first steps in public health were made: attempts to cope with the unsanitary conditions of the cities and, by means of quarantine, to limit the spread of disease; the establishment of hospitals; and provision of medical care and social assistance.

The Renaissance. Centuries of technological advance culminated in the 16th and 17th centuries in a number of scientific accomplishments. Educated leaders of the time recognized that the political and economic strength of the state required that the population maintain good health. No national health policies were developed in England or on the Continent, however, because the government lacked the knowledge and administrative machinery to

Metro-
politan
planning
council

Concerns

Greek
theories
of disease
causation

The Black
Death

carry out such policies. As a result, public health problems continued to be handled on a local community basis, as they had been in medieval times.

Scientific advances of the 16th and 17th centuries laid the foundations of anatomy and physiology. Observation and classification made possible the more precise recognition of diseases. The idea that microscopic organisms might cause communicable diseases had begun to take shape.

Among the early pioneers in public health medicine was John Graunt, who in 1662 published a book of statistics, which had been compiled by parish and municipal councils, that gave numbers for deaths and sometimes suggested their causes. Inevitably the numbers were inaccurate but a start was made in epidemiology.

National developments in the 18th and 19th centuries. Nineteenth-century movements to improve sanitation occurred simultaneously in several European countries and were built upon foundations laid in the period between 1750 and 1830. From about 1750 the population of Europe increased rapidly, and with this increase came a heightened awareness of the large numbers of infant deaths and of the unsavoury conditions in prisons and in mental institutions.

This period also witnessed the beginning and the rapid growth of hospitals. Hospitals founded in Britain, as the result of voluntary efforts by private citizens, helped to create a pattern that was to become familiar in public health services. First, a social evil is recognized and studies are undertaken through individual initiative. These efforts mold public opinion and attract governmental attention. Finally, such agitation leads to governmental action.

This era was also characterized by efforts to educate people in health matters. In 1852 Sir John Pringle published a book that discussed ventilation in barracks and the provision of latrines. Two years earlier he had written about jail fever (now thought to be typhus), and again he emphasized the same needs as well as personal hygiene. In 1754 James Lind published a treatise on scurvy, a disease caused by a lack of vitamin C.

As the Industrial Revolution developed, the health and welfare of the workers deteriorated. In England, where the Industrial Revolution and its bad effects on health were first experienced, there arose in the 19th century a movement toward sanitary reform that finally led to the establishment of public health institutions. Between 1801 and 1841 the population of London doubled; that of Leeds nearly tripled. With such growth there also came rising death rates. Between 1831 and 1844 the death rate per thousand increased in Birmingham from 14.6 to 27.2; in Bristol, from 16.9 to 31; and in Liverpool, from 21 to 34.8. These figures were the result of an increase in the urban population that far exceeded available housing and of the subsequent development of conditions that led to widespread disease and poor health.

Around the beginning of the 19th century humanitarians and philanthropists in England worked to educate the population and the government on problems associated with population growth, poverty, and epidemics. Thomas Malthus wrote in 1798 about population growth, its dependence on food supply, and the control of breeding by contraceptive methods. The utilitarian philosopher Jeremy Bentham propounded the idea of the greatest good of the greatest number as a yardstick against which the morality of certain actions might be judged. Thomas Southwood Smith founded the Health of Towns Association in 1839, and by 1848 he served as a member of the new government department, then called the General Board of Health. He published reports on quarantine, cholera, yellow fever, and the benefits of sanitary improvements.

The Poor Law Commission, created in 1834, explored problems of community health and suggested means for solving them. Its report, in 1838, argued that "the expenditures necessary to the adoption and maintenance of measures of prevention would ultimately amount to less than the cost of the disease now constantly engendered." Sanitary surveys proved that a relationship exists between communicable disease and filth in the environment, and it was said that safeguarding public health is the province of the engineer rather than of the physician.

The Public Health Act of 1848 established a General Board of Health to furnish guidance and aid in sanitary matters to local authorities, whose earlier efforts had been impeded by lack of a central authority. The board had authority to establish local boards of health and to investigate sanitary conditions in particular districts. Since this time several public health acts have been passed to regulate sewage and refuse disposal, the housing of animals, the water supply, prevention and control of disease, registration and inspection of private nursing homes and hospitals, the notification of births, and the provision of maternity and child welfare services.

Advances in public health in England had a strong influence in the United States, where one of the basic problems, as in England, was the need to create effective administrative mechanisms for the supervision and regulation of community health. In America recurrent epidemics of yellow fever, cholera, smallpox, typhoid, and typhus made the need for effective public health administration a matter of urgency. The so-called Shattuck report, published in 1850 by the Massachusetts Sanitary Commission, reviewed the serious health problems and grossly unsatisfactory living conditions in Boston. Its recommendations included an outline for a sound public health organization based on a state health department and local boards of health in each town. In New York City (in 1866) such an organization was created for the first time in the United States.

Nineteenth-century developments in Germany and France pointed the way for future public health action. France was preeminent in the areas of political and social theory. As a result the public health movement in France was deeply influenced by a spirit of public reform. The French contributed significantly to the application of scientific methods for the identification, treatment, and control of communicable disease.

Although many public health trends in Germany resembled those of England and France, the absence of a centralized government until after the Franco-German War did cause significant differences. After the end of that war and the formation of the Second Reich, a centralized public health unit was formed. Another development was the emergence of hygiene as an experimental laboratory science. In 1865 the creation at Munich of the first chair in experimental hygiene signaled the entrance of science into the field of public health.

There were other advances. The use of statistical analysis in handling health problems emerged. The forerunner of the United States Public Health Service came into being, in 1878, with the establishment in the United States of port quarantine on a national basis and with assignment of enforcement of the quarantine to the Surgeon General of the Marine Hospital Service. (Port quarantine was the isolation of a ship at port for a limited period to allow time for the manifestation of disease.)

Developments from 1875. The work of an Italian bacteriologist, Agostino Bassi, with silkworm infections early in the 19th century prepared the way for the later demonstration that specific organisms cause a number of diseases. Some questions, however, were still unanswered. These included problems related to variations in transmissibility of organisms and in susceptibility of individuals to disease. Light was thrown on these questions by discoveries of human and animal carriers of infectious diseases.

In the last decades of the 19th century the French chemist Louis Pasteur, the Germans Ferdinand Julius Cohn and Robert Koch, and others developed methods for isolating and characterizing bacteria; the English surgeon Joseph Lister developed concepts of antiseptic surgery; the English physician Ronald Ross identified the mosquito as the carrier of malaria; a French epidemiologist, Paul-Louis Simond, provided evidence that plague is primarily a disease of rats spread by rat fleas; and two Americans, Walter Reed and James Carroll, demonstrated that yellow fever is caused by a filterable virus carried by mosquitoes. Thus, modern public health and preventive medicine owe much to the early medical entomologists and bacteriologists. A further debt is owed bacteriology because of its offshoot, immunology.

In 1881 Pasteur established the principle of protective

Developments in Germany and France in the 19th century

The health of industrial workers

19th-century developments in bacteriology

vaccines and thus stimulated an interest in the mechanisms of immunity. The development of microbiology and immunology had immense consequences for community health. In the 19th century the efforts of health departments to control contagious disease consisted in attempts to improve environmental conditions. As bacteriologists identified the microorganisms that cause specific diseases, progress was made toward the rational control of specific infectious diseases.

In the United States the diagnostic bacteriologic laboratory was developed—a practical application of the theory of bacteriology, which evolved largely in Europe. These laboratories, established in many cities to protect and improve the health of the community, were a practical outgrowth of the study of microorganisms, just as the establishment of health departments was an outgrowth of an earlier movement toward sanitary reform. And just as the health department was the administrative mechanism for dealing with community health problems, the public health laboratory was the tool for the implementation of the public health program. Evidence of the effectiveness of this new phase of public health may be seen in statistics of immunization against diphtheria—in New York City the mortality rate due to diphtheria fell from 785 per 100,000 in 1894 to 1.1 per 100,000 in 1940.

While improvements in environmental sanitation during the first decade of the 20th century were valuable in dealing with some problems, they were of only limited usefulness in solving the many health problems found among the poor. In the slums of England and the United States malnutrition, venereal disease, alcoholism, and other diseases were widespread. Nineteenth-century economic liberalism held that increased production of goods would eventually bring an end to scarcity, poverty, and suffering. By the turn of the century, it seemed clear that deliberate and positive intervention by reform-minded groups, including the state, also would be necessary. For this reason many physicians, clergymen, social workers, public-spirited citizens, and government officials promoted social action. Organized efforts were undertaken to prevent tuberculosis, lessen occupational hazards, and improve children's health.

The first half of the 20th century saw further advances in community health care, particularly in the welfare of mothers and children and the health of schoolchildren, the emergence of the public health nurse, and the development of voluntary health agencies, health education programs, and occupational health programs.

In the second half of the 19th century two significant attempts were made to provide medical care for large populations. One was by Russia, and took the form of a system of medical services in rural districts; after the Communist Revolution, this was expanded to include complete government-supported medical and public health services for everyone. Similar programs have since been adopted by a number of European and Asian countries. The other attempt was prepayment for medical care, a form of social insurance first adopted toward the close of the 19th century in Germany, where prepayment for medical care had long been familiar. A number of other European countries adopted similar insurance programs.

In the United Kingdom, a royal-commission examination of the Poor Law in 1909 led to a proposal for a unified state medical service. This service was the forerunner of the 1946 National Health Service Act, which represented an attempt by a modern industrialized country to provide services to all people.

In recent years prenatal care has made a substantial contribution to preventive medicine, for it is hoped that through the education of mothers the physical and psychological health of families may be influenced and passed on to succeeding generations. Prenatal care provides the opportunity to educate the mother in personal hygiene, diet, exercise, the damaging effects of smoking, the careful use of alcohol, and the dangers of drug abuse.

Public health interests also have turned to such disorders as cancer, cardiac disease, thrombosis, lung disease, and arthritis, among others. There is increasing evidence that several of these disorders are caused by factors in

the environment; for example, the association of cigarette smoking with certain lung and cardiovascular diseases. Theoretically, they are preventable if the environment can be altered. Health education is of great importance and is a responsibility of national and local government agencies as well as voluntary bodies. Life expectancy has increased in almost every country, except where public health standards are low.

MODERN ORGANIZATIONAL AND ADMINISTRATIVE PATTERNS

International organizations. Since ancient times, the spread of epidemic disease demonstrated the need for international cooperation for health protection. Early efforts toward international control of disease appeared in national quarantines in Europe and the Middle East. The first formal international health conference, held in Paris in 1851, was followed by a series of similar conferences aimed at drafting international quarantine regulations. A permanent international health organization was established in Paris in 1907 to receive notification of serious communicable diseases from participating nations, to transmit this information to the member nations, and to study and develop sanitary conventions and quarantine regulations on shipping and train travel. This organization was ultimately absorbed by the World Health Organization (WHO) in 1948.

In the Americas, the organization of international health probably began with a regional health conference in Rio de Janeiro in 1887. From 1889 onward there were several conferences of American countries, which led ultimately to the establishment of the Pan-American Sanitary Bureau; this was made a regional office of WHO in 1949, when it became known as the Pan-American Health Organization.

The rise and decline of health organizations has been influenced by wars and their aftermaths. After World War I, a Health Section of the League of Nations was established and functioned until World War II. After the war, the United Nations Relief and Rehabilitation Administration (UNRRA) was set up; it processed displaced persons in such a way as to prevent the spread of disease. It was responsible for the planning steps that led to the establishment in 1948 of the World Health Organization as a special agency of the United Nations. WHO is concerned with physical, mental, and social well-being and not merely with the absence of disease.

The work of WHO is carried out under the direction of the World Health Assembly, which has representatives from the member states. The first assembly gave consideration to diseases and problems that exist in large areas of the world and that lend themselves to international action. Malaria, tuberculosis, venereal disease, the promotion of health, environmental conditions responsible for a significant proportion of deaths, and nutrition were given priority. Other areas of need have been included since.

Among important functions of the organization are the advisory services offered to governments through its regional staff. Regional offices in a number of countries, both industrialized and developing, as well as local representatives in many developing countries, help WHO maintain contact with needs and sources of financial aid. In specialized fields, a number of expert committees consider specific questions.

WHO maintains close relationships with other United Nations agencies, particularly the United Nations Children's Fund (UNICEF) and the Food and Agriculture Organization (FAO), and with international labour organizations. From its inception in 1946, UNICEF focused its aid on maternal and child health services and the control of infections, especially in children. Priority has been given to the production of vaccines, the institution of environmental sanitation, the provision of clean water, and the training of local personnel in their own countries (especially in rural areas). Aid is channeled through organized health services in developing countries. Recent efforts have concentrated on persuading governments to undertake national surveys to identify the basic needs of their children and to devise appropriate national policies.

The work of WHO includes three main categories of

Poverty
and
disease

Establish-
ment
of WHO

activities. First, it is a clearinghouse for information about disease throughout the world, and it has developed a uniform system for reporting diseases and causes of death. It has established internationally accepted standards for drugs and drawn up a list of "essential" (effective, cheap, and reliable) drugs. It has sponsored and financed many research projects throughout the world. Second, WHO has promoted mass campaigns to control epidemic and endemic diseases, a substantial number of which have been quite successful. Third, WHO attempts to strengthen and expand the public health administration and services of member nations by providing technical advice, teams of experts to carry out surveys and demonstrate projects, and aid in support of regional and national health development projects.

Developed nations. Methods of health administration vary from country to country. Major health functions are frequently grouped in a department that is responsible for health and for related functions. In the United Kingdom they are carried out by the Department of Health and Social Security; in the United States the Department of Health and Human Services controls the programs covered by national legislation.

Few central departments of health are all-embracing; other departments also operate medical programs of some sort. No country places the health services of its military forces under the central health agency. Because unity of control at the centre is impracticable, coordination is important. Central administration is further complicated in federal systems. In the United States there are 50 states, no two of which have the same patterns of health organization.

Patterns shared. The official responsible for the administration of national health affairs is in most cases a member of the Cabinet. Advisory councils are frequently used to bring the ideas of leading scientists, health experts, and community leaders to bear on major national health problems.

An organization that provides basic community health services under the direction of a medical officer is called a local health unit. It is usually governed by a local authority. Its programs may include maternal and child health, communicable-disease control, environmental sanitation, maintenance of records for statistical purposes, health education of the public, public health nursing, medical care, and, often, school health services. The local health unit can provide the administrative framework for a wider range of community health services, including the care of the aged, of the physically handicapped, and of the chronically ill and mental health services. Although social welfare services may be provided by a separate agency, there are advantages in amalgamating health and welfare services, because a family's health and social problems tend to be interrelated. In England welfare and public health are often integrated at the local level, whereas in the United States they are almost always separate.

The population served by a local health unit may be only a few thousand or several hundred thousand. There are substantially different problems involved in administering health services for a large rural area that is sparsely populated and a municipality with a population of one or two million.

One problem of administering local health services is the question of whether they should be run by independent local authorities or organized regionally to ensure coordination and effective referral and to avoid duplication of services.

Medical care is provided as a public service to some degree in most countries. It may be limited to the hospitalization of persons afflicted with certain ailments—for example, mental disease, tuberculosis, chronic illness, and acute infections. Comprehensive health services may be provided for some specific population groups, as in Canada and the United States, where the federal government provides care for Indians and Eskimos. Many countries have compulsory medical insurance, and some combine the socialization of hospitals with medical insurance covering general medical care, as in Denmark. Full-scale socialization of health services exists in a few countries,

including the United Kingdom and New Zealand.

In countries such as The Netherlands and the United States, where voluntary and nonprofit organizations support a considerable share of the health services and operate most of the general hospitals, there is pluralism in health administration. This makes coordination difficult, but voluntary effort has the advantages of involving citizens directly in the development of health services and of promoting experimentation in administration.

There is a trend toward regional planning of comprehensive health services for defined populations. In an idealized plan, the first level of contact between the population and the system, which can be called primary care, is provided by health personnel who work in community health centres and who reach beyond the health centres into the communities and homes with preventive, promotive, and educational services. At the next level of care, specialists in community hospitals provide secondary care for patients referred from the primary-care centres. Finally, tertiary, or superspecialty, care is provided by a major medical centre. The various levels of this regional scheme are linked by a two-way flow of medical records, patients, and health personnel. Regionalization has been most fully achieved in Europe and least so in North America, where voluntary hospitals provide most of the short-term general services and retain autonomy in their administration.

Variations. Among the developed nations, there is substantial variation in the organization and administration of health services. The United Kingdom, for example, has a National Health Service with substantial autonomy given to local government for implementation. The United States has a pluralistic approach to health services, in which local, state, and national governments have varying areas of responsibility, with the private sector playing a prominent role.

During the first half of the 20th century in Britain, the emphasis shifted gradually from environmental toward personal public health. A succession of statutes, of which the Maternity and Child Welfare Act (1918) was probably the most important, placed responsibility for most of the work on county governments. National health insurance (1911) gave benefits to 16,000,000 workers and marked the beginning of a process upon which the National Health Service Act (1946) was built.

The National Health Service Act provided comprehensive coverage for most of the health services, including hospitals, general practice, and public health. The service remained at the periphery, however, in three types of care: (1) Primary medical care is given by family physicians or general practitioners. This service is organized locally by an executive council. Each general practitioner is responsible for providing primary care to a group of people on a particular registry. (2) Specialist consultation and outpatient and inpatient treatment are provided in hospitals under the direction of regional authorities. A later concept makes each district general hospital responsible for providing hospital services for a defined population. (3) Services, such as health visiting, home nursing, home helps, domiciliary midwifery, the prevention of illness, and the provision of health centres are the responsibility of local authorities.

In the former Soviet Union the protection and promotion of public health was the responsibility of the state. There was free public access to all forms of medical care. The principles of the health services were complete integration of curative and preventive services, medicine as a social service, preventive programs, health centres or polyclinics (clinics in which a variety of diseases were handled), and community participation.

The public health services for the Soviet Union were directed by the Ministry of Health. Each of 15 republics of the union had its own ministry. Each republic was divided into *oblasti* (provinces), which in turn were divided into *rayony* (municipalities) and finally into *uchastoki* (districts). Each subdivision had its own health department accountable to the next highest division (see above *The practice of public health*).

There were well-established referral procedures, from the polyclinics and smaller hospitals in the *uchastoki* to the

Levels of health care

Departments of health

Types of public medical service

larger *rayon* hospitals, and from feldshers (paramedical personnel trained in medical care) and other paramedical personnel to internists and pediatricians and, when necessary, to more highly specialized personnel.

The health services of the United States can be considered at three levels: local, state, and federal.

Levels of service

Locally in cities or counties, there is substantial autonomy within broad guidelines developed by the state. The size and scope of local programs vary, but some of their functions are control of communicable diseases; clinics for mothers and children, particularly for certain preventive and diagnostic services; public health nursing services; environmental health services; health education; vital statistics; community health centres, hospitals, and other medical care facilities; community health planning and coordination.

At the state level, a department of health is charged with overall responsibility for health, though a number of agencies may actually be involved. The state department of health usually has five functions: public health and preventive programs; medical and custodial care such as the operation of hospitals for mental illness; expansion and improvement of hospitals, medical facilities, and health centres; licensure for health purposes of individuals, agencies, and enterprises serving the public; and financial and technical assistance to local governments for conducting health programs.

At the federal, or national, level, the Public Health Service of the Department of Health and Human Services is the principal health agency, but several other departments have health interests and responsibilities. Federal health agencies accept responsibility for improving state and local services, for controlling interstate health hazards, and for working with other countries on international health matters. The federal government also has the following specific responsibilities: (1) protecting the United States from communicable diseases from abroad; (2) providing for the medical needs of military personnel, veterans, merchant seamen, and American Indians; (3) protecting consumers against impure or misbranded foods, drugs, and cosmetics; and (4) regulating production of biological products, such as vaccines. In addition, the federal government promotes and supports medical research, health services, and educational programs throughout the country.

Voluntary effort is a significant part of health work in the United States. There are more than 100,000 voluntary agencies in the health field functioning mostly at the local level but also at state and national levels. Supported largely through private sources, these agencies contribute to programs related to education, research, and health services.

Channels for medical care

Medical care is provided and paid for through many channels, including public institutions, such as municipal, county, state, and federal health centres, hospitals, and medical care programs, and through private hospitals and private practitioners working either alone or, increasingly, in groups. Generally, medical care is financed by public funds, voluntary health insurance, or personal payment. There is a trend away from the traditional fee-for-service payment to individual practitioners toward prepaid-care systems including health teams working at community, health centre, and hospital levels.

Thus, in the United States there is great variety in the content, scope, and quality of health services. These services are provided by several independent agencies. In effect, however, they constitute a working partnership for the protection and promotion of human health.

Recently, two factors have contributed to rapid change in the orientation of health services in the United States. One of these is an increasing awareness that, while the existing system of health services provides high quality care for many, there are others for whom the care is either lacking or unsatisfactory. The second factor is that of steeply rising costs of medical care. These two issues have led to reconsideration of the entire system of personal medical care and proposals for new systems of providing and financing health care.

Developing nations. *Patterns shared.* Developing countries have sometimes been influenced in their approaches to health care problems by the developed

countries that have had a role in their history. The countries in Africa and Asia that were once colonies of Britain have educational programs and health-care systems that reflect British patterns, though there have been adaptations to local needs. Similar effects may be observed in countries influenced by France, The Netherlands, and Belgium.

Despite variations from country to country, a common, if somewhat idealized, administrative pattern may be drawn for developing countries. All health services, except for a small amount of private practice, are under a ministry of health, in which there are about five bureaus, or departments—hospital services, health services, education and training, personnel, and research and planning. Hospital and health services are distributed throughout the country. At the periphery of the system are dispensaries, or health outposts, often manned by one or two persons with limited training. The dispensaries are often of limited effectiveness and are upgraded to full health centres when possible. Health centres and their activities are the foundation of the system. Health centres are usually staffed by auxiliaries who have four to 10 years of basic education plus one to four years of technical training. The staff may include a midwife, an auxiliary nurse, a sanitarian, and a medical assistant. The assistants, trained in the diagnosis and treatment of sickness, refer to a physician the problems that are beyond their own competence. Together, these auxiliaries provide comprehensive care for a population of 10,000 to 25,000. Several health centres together with a district hospital serve a district of about 100,000 to 200,000 people. All health services are under the responsibility of the district medical officer, who, assisted by other professional and auxiliary personnel, integrates the health efforts into a comprehensive program.

The typical administrative pattern

Of central importance is the distribution of responsibilities between auxiliaries and professionals. The auxiliaries, by handling the large number of relatively simple problems, allow the professionals to look after only the more complex problems, to supervise and teach the auxiliaries, and to plan and manage the programs.

The district hospital is dependent on a regional hospital, to which patients with complex problems can be referred for more specialized services. Administrative direction of both regional health services and regional hospital services can be combined at this level under a regional medical officer. The central administration of the ministry of health provides policies and guidance for an entire health service and, in some instances, also provides a central planning unit.

Problems of transportation and communication over great distances, shortages of staff and other resources, and inadequacies in staff preparation and motivation often lead to malfunctions in the system. Nonetheless, the public health services developed in African and Asian countries have generally provided a sound basis for future development within the framework of national development.

Variations. The organization of public health services in Latin-American countries differs substantially from those of Africa and Asia; these differences are an expression of their different historical backgrounds. The Latin-American countries are generally more affluent than those of Asia and Africa. Private practice is more widespread, and private or voluntary agencies are more prominent. Health services are provided largely by local and national governments. Many Latin-American countries also have systems of clinics and hospitals for workers financed by employers and workers. The distribution of health services, with health centres, hospitals, and preventive services, is roughly similar to Africa and Asia. The Latin-American countries, however, have used auxiliaries less than African and Asian countries. Latin America has pioneered in the development of health-planning methods. Chile has one of the most advanced approaches to health planning in the world.

Differences between Latin America and Africa

Thailand was never colonized and therefore has no historical influence favouring any particular pattern of health services. The Thai Ministry of Health has a well-developed system of hospitals and health centres across the country to serve both rural and urban people. It differs from the pattern described in the previous section in that, despite

the extreme shortages of physicians and nurses in rural areas, the nation has been reluctant to use auxiliaries for medical care. It does, however, use auxiliary midwives and sanitarians. Hospital services and public health services have separate administration. Within the public health services, there are a number of separate divisions—e.g., for tuberculosis, sexually transmitted diseases, and nutrition—each with its own staff, budget, and facilities. The trend elsewhere has been away from relatively independent, disease-oriented approaches and toward integrated systems in which the same network of health services handles most problems.

Health problems and obstacles. The difficulties of providing health services for the people of the developing nations involve a cluster of interrelated problems. These arise from the nature of the diseases and hazards to health, insufficient and maldistributed resources, the design of health service systems, and the education of health personnel in those systems. Woven through the health programs of the developing nations and complicating them at both family and national levels are the pressures associated with rapidly growing populations.

There are differences not only in the kinds of diseases of different countries but also in the rates at which they occur and in the age groups involved. Life expectancy in some countries is less than half that in others, principally because of high death rates among small children in the developing countries. In much of Southeast Asia, for example, 40 percent of children die by their fourth year, a death rate not reached until age 60 in North America. The infant (under one year of age) mortality rate in Central and South America is two to four times that in North America, and the death rate in children one to four years of age is as much as 25 times greater. The differences for Central Africa are even more striking: infant mortality in some areas has been 12 times that in the United States, and the mortality in preschool children has been more than 60 times the U.S. figure.

The principal causes of sickness and death among small children in the developing world are diarrhea, respiratory infections, and malnutrition, all of which are intimately related to culture, custom, and economic status. Malnutrition may result from food customs when taboos and simple oversight lead to deprivation of children. Gastroenteritis (inflammation of the lining of the stomach and intestines, usually with accompanying diarrhea) and respiratory infections are often due to infectious organisms that are not susceptible to antibiotics. The interrelationships of these diseases increase the complexity of treating them. Malnutrition is often the underlying culprit; not only does it cause damage itself, such as retardation of physical and mental development, but it also seems to set the stage for other illnesses. A malnourished child develops gastroenteritis, inability to eat, further weakness, and then dehydration. The weakened child is susceptible to a lethal infection, such as pneumonia. Or, to complete the vicious circle, infection can affect protein metabolism in ways that contribute to malnutrition.

Another factor that contributes to this is family size. Malnutrition, with associated death and disability, occurs most often in children born into large and poorly spaced families. The resulting high death rate among small children often reinforces the tendency of parents to have more children. People are not inclined to limit the size of their families until it is apparent that their children have a reasonable chance of survival. Thus, there is a fertility-mortality cycle in which high fertility, reflected in large numbers of small children crowded into a poor home, leads to high childhood mortality, which, in turn, encourages high fertility. This is the basis of the belief that population-control programs should include effective means of reducing unnecessary deaths among children.

Among limitations of resources, shortages of trained personnel are among the most important; ratios of population to physicians, nurses, and beds provide an indication of the seriousness of these deficiencies and also of the great differences from country to country. Thus, the proportion of population to physicians in developing countries varies drastically.

Money is a crucial factor in health care—it determines how many health personnel can be trained, how many can be maintained in the field, and the resources that they will have to work with when they are there. Governmental expenditures on health care vary greatly from country to country.

As it attempts to provide health care for its people, a nation, on the one hand, must meet the urgent and complex problems, such as obstetric and surgical emergencies for which hospital care is essential. On the other hand, it must reach into the communities and homes to find those who need care but do not seek it and must discover the causes of such diseases as malnutrition and gastroenteritis.

Education of health personnel. In the education of health personnel, a particular set of problems emerges. Educational programs for auxiliaries are suited to the local situation, perhaps because they were not established in the more developed nations. Medical and nursing education, on the other hand, is similar to that of the more advanced countries, and it prepares students better for working in industrialized nations than in their own. This misfit between education and the jobs to be done has probably contributed substantially both to the ineffectiveness of health service systems and to the migration of professional personnel to the more developed countries.

PROGRESS IN PUBLIC HEALTH

Developed nations. Among the more developed nations the following trends are apparent.

Increasing interest of national governments. Formerly, governments were chiefly concerned with basic health problems, such as environmental sanitation, medical care of the poor, quarantine, and the control of communicable diseases. Gradually, they have extended their activities into the field of medical care services in the home, clinic, and hospital, so as to provide comprehensive health care for entire communities. Three factors have influenced this trend: (1) the nongovernmental voluntary agencies have been unable to meet the rising cost of medical care; (2) there is an increasing appreciation of the economic loss to a country from sickness; and (3) there is an increasing public interest in social services. Health and social welfare are now recognized as complementary, and social legislation tends to cover both areas. There is an administrative trend toward a close cooperation between health and social welfare services.

Changing concepts of preventable disease. Until recently, the term preventable disease referred to a circumscribed group of infectious diseases. The term has acquired a broader meaning, however, as epidemiological methods are applied to other conditions. Preventive health services now deal with a wide range of health hazards, such as malignant tumours, rheumatism, cardiovascular diseases, other chronic and degenerative diseases, and even accidents.

Integration of preventive and medical care services. Medical care had its origin in the humanitarian motive of caring for the sick, while preventive health services sprang from the need to protect a healthy environment from epidemic diseases. They grew apart, but recently the trend has been to integrate them within a comprehensive health service. Such an integration was the fundamental principle of public health in the U.S.S.R., in which all local health services were centred in the district hospital under one administration. In European countries, especially in rural areas, the two branches are brought together by the local medical practitioner. The focal point of many discussions on medical care is the role that the hospital should play in health services. Many feel that its influence at present is too restricted and that it should spread beyond its walls to health centres and homes.

Provisions directed toward better mental health. Mental health now has a place in the preventive services. Improvements in arrangements for mental health include the provision of outpatient clinics and inpatient accommodations at general hospitals for early mental cases, an increase in child-guidance and marriage-guidance clinics, and schemes for the care of alcoholics and drug addicts. There have also been significant developments in the treat-

Infant
mortality
rates

Ratios of
medical
personnel
to
population

Factors
in the
increase
of govern-
ment
involvement

Compre-
hensive
health
service

ment of maladjusted members of society. Gains in understanding of psychoneuroses by general practitioners and the development of research facilities are also noteworthy.

Growing emphasis on health education. Many countries have expanded their commitment to health education, usually in cooperation with voluntary agencies. The most effective work is carried out at the local level, especially in schools. The trend is toward an expansion of health education as an essential preventive health service.

The biostatistical, epidemiological approach. A statistical service is essential in planning, administering, and evaluating health services. The interest of public authorities in medical-care schemes has increased the importance of statistics on the incidence of diseases and other problems, as well as the epidemiology necessary to combat them. Both are vital in the planning, organization, and evaluation of medical-care schemes. Traditionally, the epidemiological method was used for infectious diseases, but it is now being used increasingly for noninfectious diseases and the problems of medical care.

Changes resulting from an aging population. In more affluent nations, an increase in older age groups brings about the need for public health facilities to provide special services for them. Health care of the elderly includes measures to prevent premature aging and the chronic and degenerative diseases and to confront the psychological problems resulting from loneliness and inactivity. Geriatric clinics have been set up to meet these needs and to conduct research into the process of senescence.

Concern regarding the quality of the environment. There is widespread concern about environmental deterioration. Controlled atomic radiation has created new hazards to health, such as the potential pollution of air or water by radioactive discharges, the possible effects from radioactive fallout on the public generally, and the dangers to workers in atomic installations in industry. A growing population requires an increase in industrial and commercial activities, which add to the volume of pollutants that threaten the atmosphere, rivers, lakes, and oceans and have destructive effects on natural ecology. Individual countries have taken steps toward the control of environmental deterioration, and means of international regulation have also been proposed.

Developing nations. In view of the large numbers of serious health problems in the developing nations and their limited resources for dealing with them, it is understandable that along with substantial progress there would be some stagnation, or even regression.

Communicable-disease control. Smallpox and malaria are examples of diseases that have been brought under closer control throughout the world. For other diseases, such as hepatitis (liver inflammation), rabies, leprosy, and sleeping sickness, there have been important growths in understanding that may contribute to their eventual control.

Disease problems that await solution. El Tor cholera, which has appeared in epidemic form in previously uninvolved areas, represents one of the most serious challenges to public health. Venereal disease, an old problem, has increased in incidence. Certain parasitic diseases have spread as humans have brought about changes in their environment—the increase in schistosomiasis (infestation with blood fluke by means of snails as the intermediate hosts) in irrigation and man-made lake areas is an example. Widespread malnutrition, particularly protein-calorie malnutrition in small children, remains a problem. Protein-rich food supplements and more effective educational programs are being developed to combat it.

Family health. The problems of rapidly growing populations have important consequences at both the family and the national level. Problems of maternal and child health, human reproduction, and human genetics, including family planning, are now seen as aspects of the greater problem of the health of the whole family as a single and fundamental social unit. Accordingly, family health is a matter deserving high priority among the public health services.

Health manpower. There is widespread recognition of inadequacies in both number and education of health personnel. The trend is toward coordinating the education of

health personnel with the particular health service in which they will function. This trend requires close relationships between educational institutions and the agencies responsible for health services.

Comprehensive community health services. The fragmentation of earlier health service organizations, such as programs concerned with only a single disease and the separation of curative and preventive services, is giving way to more comprehensive organizational patterns. Health promotion, disease prevention, and the curing and rehabilitation of the ill are brought together into one network of integrated services that reaches to the community level.

National health planning. Complex decision making is involved in allocating limited health service resources to large numbers of people. As a result, there has been an increasing emphasis on the health-planning process and on the design of more effective health-service systems. A number of countries have established health-planning units in the ministry of health or the national planning organization. An important aspect of national health planning is the close coordination between planning, budgeting, implementing, and evaluating programs.

(J.H.Br./P.Rh.)

Clinics

A clinic is an organized medical service offering diagnostic, therapeutic, or preventive treatment to ambulatory patients. Often in Europe and occasionally in the United States the term covers the entire teaching centre, including the hospital and the ambulatory-patient facilities. The medical care offered by a clinic may or may not be connected with a hospital.

The term clinic may be used to designate all the activities of a general clinic or only a particular division of the work; e.g., the psychiatric clinic, neurology clinic, or surgery clinic. The entire activity when connected with a hospital is called the outpatient department, and the specific subdivisions are referred to as clinics.

The first clinic in the English-speaking world, the London Dispensary, was founded in 1696 as a central means of dispensing medicines to the sick poor whom the physicians were treating in the patients' homes. The New York City, Philadelphia, and Boston dispensaries, founded in 1771, 1786, and 1796, respectively, had the same objective. Later, for the sake of convenience, physicians began to treat their free patients at the dispensary. The number of such clinics did not increase rapidly, and as late as 1890 only 132 were operating in the United States. The impetus for the mushroomlike growth that has occurred since that time came with the rapid growth of hospitals and also from the public health movement.

HOSPITAL CLINICS

During the late 1800s the modern concept of a hospital began to take shape. During this period some of the hospitals connected with medical schools inaugurated outpatient departments for the purposes of teaching and charity. The advantages of providing ambulatory care close to the facilities of a hospital became apparent, and such hospital clinics multiplied rapidly.

The organization of a hospital clinic in general follows that of the inpatient facilities. Hospital clinics are primarily concerned with acute diseases, and the physicians in the clinics are usually the same physicians who treat inpatients in the hospital.

In many hospital clinics, especially those in countries that do not have national health insurance programs, care is made available only to the medically indigent, and no professional fee is charged. Practically all such clinics, however, charge a small registration fee if the patient is financially able to pay; income from such fees helps pay operating costs. A number of successful attempts have been made to extend hospital clinic care to paying patients. Most of this effort has been in the area of lower income groups although in a few hospitals no limit is placed on income in determining eligibility for care. The hospitals of the University of Chicago, for example, began operating a clinic on such a basis in 1928.

New hazards to health

Population growth

The London Dispensary

PUBLIC HEALTH CLINICS

The first clinics

The public health movement has been mainly concerned with preventive medicine, child and maternal health, and other medical problems affecting broad segments of the population. The first public health clinics were established in the late 19th century. In 1890 A. Pinard set up a maternal dispensary or antenatal clinic at the Maternité Baudelocque in Paris. Milk distribution centres were set up in France by J. Comby (1890) and in Britain by F.D. Harris (1899). Infant welfare clinics were established in Barcelona (1890); and clinics for older children were founded in St. Pancras, London, by J.F.J. Sykes (1907).

Unlike hospital clinics, which have had their greatest growth in the cities, public health clinics are located chiefly in smaller towns and villages. In the United States the first great movement in creating public health clinics resulted in the founding of the National Association for the Study and Prevention of Tuberculosis in 1904. It was the association's goal to study and prevent tuberculosis by making clinic facilities available for free examination and treatment in every city and county. Other nationwide, private health agencies in specialized medical fields quickly adopted this method to improve the quality and extend the quantity of medical service in their fields. Local governmental health units operate similar clinics for the prevention of communicable disease and long-term illness. Such clinics are generally concerned with one particular medical interest, usually one of the following: tuberculosis, sexually transmitted diseases, prenatal care, well-baby care, teeth, tonsils, eyes, crippled children, and mental health. There is a tendency toward the establishment of traveling clinics, such as dental clinics for schoolchildren. Often no charge is made for service in public health clinics, and for many medical conditions no income restrictions are imposed. A few are operated in connection with hospitals, but most such clinics use public buildings or space furnished by welfare and other social agencies. Financial support is received mostly from the same sources.

Traveling clinics

PRIVATE CLINICS

The advantages of group medical service, with facilities and technical personnel beyond the means of an individual practitioner plus the benefit of group consultation, have encouraged the establishment of pay or private clinics. Such a clinic is essentially a voluntary association of physicians engaged in the practice of medicine on an organized group basis. Common administration and facilities are used, and the resulting expense and income are shared according to a predetermined plan. To be classified as a group clinic the relationship between each physician and the organization must be defined in a legal agreement. The relationship usually takes the form of a partnership. Several of these, such as the Mayo Clinic, in Rochester, Minn., have achieved a national reputation and attract patients from a wide area. Most of these organized group clinics are general clinics—*i.e.*, they have several of the different medical specialties represented on their staffs. A number of private clinics, however, limit their work to one medical specialty. An enterprise of special interest is the London Clinic, established in 1936 by a group of prominent consultant surgeons and physicians who wished to make available to their private patients a place where the comforts and privacy of a nursing home could be combined with facilities for diagnosis and therapy such as exist only in the larger general hospitals.

Usually the group is organized independently of any hospital or other agency, but in some instances such clinics own and operate their own hospital facilities. In other instances the clinic is a part of a prepaid health service plan. This latter pattern has received impetus in recent years as labour unions have set up medical clinics supported by welfare benefits contributed by employers. The United Mine Workers, for example, has established a system of such clinics in hospitals constructed by the union in the coal-mining areas of Virginia and West Virginia.

HEALTH CENTRES

In 1910 the first health centres, or multiple clinics, were established in Pittsburgh and Wilkes-Barre, Pa. Others were opened in 1913 in New York City and in 1916 in Boston and Philadelphia. In 1920 in Britain a consultative council on medical and allied services (Dawson Committee) described a health centre as an institution wherein various medical services, preventive and curative, are brought together. Under Section 21 of the National Health Service Act, 1946, local health authorities provide, equip, maintain, and staff health centres to offer facilities for all or any of the following services: general medical and dental, pharmaceutical, hospital outpatient and health education, and care of expectant and nursing mothers. The prenatal and postnatal clinics include periodic medical and dental examinations, classes in parenthood and relaxation, and welfare foods. Activities in child welfare clinics comprise education in all aspects of motherhood, periodic medical and dental examinations, advice on mental health problems, immunization and vaccination, and distribution of welfare foods.

Services

POLYCLINICS

In the U.S.S.R. the polyclinic (*poliklinika*) was created in order to combine the function of a hospital outpatient department with that of a general-practitioner health centre. Factory workers and their families could attend general polyclinics attached to major factories, and many special children's polyclinics were built in the towns. A typical urban polyclinic, which was usually associated with a hospital, included reception and waiting rooms, registration desk, and consulting and treatment rooms of the following: internist, pediatrician, gynecologist, surgeon, eye specialist, ear, nose, and throat specialist, neurologist, urologist, and dentist. Rooms were often set aside for first aid, reception of infectious cases, and women and children's welfare, as well as a dental department, drugstore, pathological laboratory, X-ray department, gymnasium, and lecture hall. There were always fairly elaborate physiotherapy departments and usually large and small operating theatres.

(R.E.Bn./J.F.Sk./Ed.)

FAMILY PLANNING CLINICS

The main purposes of the family planning service are to encourage parents to make responsible decisions about pregnancy that take into account the best interests of the family; to provide guidance to couples who wish to limit the size of their families; and to advise on the technical methods that are available for doing so. There are marked differences in attitude toward the desirability of a reduction in family size as between developed and developing countries. This difference is dominated by high infant and child mortality in most developing countries as compared with developed countries. (Ha.Sc.)

RELATED FIELDS

Nursing

The need for nursing is universal. The International Council of Nurses states that the function of nursing is four-fold—to promote health, prevent illness, restore health, and alleviate suffering—and that inherent in nursing is "the respect for life, dignity and the rights of man. It is unrestricted by considerations of nationality, race, creed, colour, age, sex, politics or social status. Nurses render

health services to individuals, the family and the community and co-ordinate their services with those of related groups." Nurses form the largest group of health service workers throughout the world. As medical science advances, the health needs of a given population change and expectations for care rise, so that the potential for nursing service increases. The World Health Organization (WHO) stated in 1986 that "nursing and midwifery personnel have an even greater potential contribution to make to health

for all than they are already making. Nurses' potential lies in their role as providers of primary health care services in community settings, clinics, schools and industries as well as in hospitals."

HISTORY OF NURSING

When Christianity developed, with its teachings of duty, love, and brotherhood, caring for the sick gained new impetus. Among the early converts to Christianity were highborn Roman matrons. Phoebe, a friend of St. Paul, "succored many"; Marcella made her palace into a monastery for women; Fabiola turned her home into a hospital. These women and others like them exemplified new intellectual freedom for women and their participation in religious and social action at a high level.

For several centuries thereafter the spread of health work moved slowly. Monastic orders grew, and from them some health services radiated; but large areas of the world were untouched, and the monastic undertaking of health work seemed to diminish. Later, military and chivalric orders were founded, supporting the Crusades and combining the making of war with charitable and hospital work. Some established companion orders for women. Monastic and chivalric orders established in various forms among health workers the beginnings of hierarchical organization, amateur participation that provided the roots of volunteerism, and a sense of calling or vocation.

Nursing can be seen to rise from two wellsprings, one scientific, the other religious and social. Acceleration of scientific advancement in health began with the 16th century. During the 19th century the "germ theory" of disease was developed; ways of treating and preventing infectious diseases, then the largest cause of death, were introduced. Anesthesia was also discovered. It is often said that research produced more medical and health knowledge in the decades after World War II than in all previous centuries combined. The growing mass of new knowledge to be applied in health services by health workers challenged the educational system for physicians, nurses, and others and strained the system of distribution of services to an awakened public.

Paralleling these scientific advances were those made through the years in social and religious action. The work of St. Vincent de Paul and, later, of John Howard typify this kind of action. After making thorough studies of conditions, they recommended marked changes in hospitals and prisons. A new dimension, planning based on facts, thus was added. St. Vincent, with St. Louise de Marillac, founded the Daughters of Charity (1633). In the following century Howard revealed deplorable conditions in nursing—filth, stealing, and ill treatment of patients. Theodor Fliedner, a 19th-century clergyman, later influenced reforms and emphasized citizens' responsibility for the health and welfare of people everywhere.

During the 19th century the movement for reform in nursing was led by Florence Nightingale, a woman of intellectual and moral power. Family contacts with humanitarian leaders and an education that included training in science, mathematics, and political economy were her preparation. She critically studied nursing as it was practiced in several countries, formulated her ideas, and wrote extensively.

In 1854 Florence Nightingale was asked by the British secretary of state at war to go to Scutari in Turkey, where absence of sewers and of laundering facilities, lack of supplies, poor food, disorganized medical service, and absence of nursing led to a death rate of more than 50 percent among wounded soldiers. Her work and that of the nurses whom she recruited brought sufficient improvement to lower the death rate to 2.2 percent. A gift of £45,000 was raised by popular subscription, and Florence Nightingale used it to establish schools of nursing at St. Thomas's Hospital in London and elsewhere.

Florence Nightingale believed nursing to be suitable as an independent career for capable, trained women, that nursing services should be administered by those with special preparation, and that relationships between physicians and nurses should be professional. She maintained that schools of nursing should be administered by nurses

with physicians as part of the hospital labour force. She believed that there was a substantial body of knowledge and skills to be learned in nursing. Nurses were to be prepared for hospital nursing and care of the sick at home, and they were to teach good health practices to patients and families.

By the end of the 19th century, the idea that a nurse needed to be educated and trained had spread to much of the Western world. In England, Scandinavia, America, and much of the British Empire, schools of nursing were generally based on training hospitals, and more nurses had become independent of religious institutions. On the Continent, however, the Motherhouse system and religious organizations often persisted. With the foundation of the Red Cross in 1863, some national Red Cross societies undertook the training of nurses and provided instruction for auxiliaries to help in time of war or emergency.

With advances in medicine, nurses followed doctors into specialties, including pediatrics, surgery, orthopedics, ophthalmology, psychiatry, and public health. The development of physicians' assistants has sometimes overlapped with the development of nurse specialists, though efforts have been made to differentiate the two categories. Pediatric nurse-practitioners, for example, now undertake functions previously performed only by physicians—taking histories, performing physical examinations on children, and running clinics for mothers and babies. There has also been experimentation with shifts in functions from medicine to nursing.

Planning is aimed at bringing nursing resources into balance with the needs for nursing personnel to implement health programs. WHO has developed a guide for such planning, published in several languages; the guide helps countries to establish goals and avoid disadvantages of haphazard growth. Coordinated planning among the several health professions and occupations is now considered to be highly desirable.

THE PRACTICE OF NURSING

Kinds of nursing. Nursing personnel and the kinds of nursing in which they engage may be classified in a variety of ways: by legal designation—registered nurses, licensed practical nurses, and, in addition, unlicensed nurses' aides; by the kind of educational preparation—vocational, technical, and professional, including graduate education—that they have received; by the kind of work they perform—institutional, community, educational, occupational, research, or journalistic; by the level of responsibility that they assume—that of staff nurse, teacher, supervisor, administrator, or consultant; and by the place of employment—hospital, physician's office, public health agency, school (school health nurses), industry, a school of nursing, including undergraduate and graduate programs. Nurses may also be classified as either generalists or specialists. Most descriptions of nursing, including the following, combine these classifications.

Hospital nursing. Hospital nursing occupies approximately two-thirds of the total nursing force in most countries. Nurses giving direct care to patients in hospitals of the United States, for example, include registered nurses, some of whom may be specialists; licensed practical nurses; and nurses' aides. Aides and sometimes practical nurses, if there are any, are called auxiliaries in many other developed countries. In large, highly organized hospitals, administrators of nursing services need special preparation and occupy demanding executive positions.

Certain areas of hospitals are highly specialized—surgical suites (operating rooms), recovery and intensive-care units, coronary-care units, and drug addiction clinics. In these, nurses giving direct patient care and their supervisors must have the additional preparation afforded by graduate study or by continuing education or staff-development (in-service-education) programs.

The rapid development of hospitals in response to advances in science and new community responsibilities requires new kinds of nursing. Some hospitals use nurse-midwives in maternity departments. Nurse specialists work in a number of other areas as well—medical, surgical, obstetric, pediatric, psychiatric, and rehabilitative; these

Methods of classifying nurses

The work of Florence Nightingale

nurses frequently cross over to other fields as consultants. For example, public health nurses are employed to ease patients' adjustment to returning home and to assure planned continuity of health care. A few nurses occupy ombudsman, or patient-advocate, positions. Some others perform research aimed at solving problems and at devising improved ways of caring for patients.

New health problems bring additional responsibilities to nurses in hospitals. Drug abuse, family planning, and terminal illness, for example, are areas that call for specialized knowledge on the part of nurses, as well as for decisions about the organization of the services in hospitals.

Outpatient
and home
care

Outpatient and home care are sometimes organized by hospitals to reduce costly hospitalization for patients. Clinics in outpatient departments may be operated almost exclusively by nurses—clinics for hypertension (high blood pressure), chronic heart disease, and poststroke and maternity patients are examples. Satellite clinics bring services close to patients' homes. These expanded services require that nurses learn assessment methods—simple physical examination and history taking, for example—understand community problems and the organization of nursing services.

Efforts to return many patients of mental hospitals to their communities have met with some success. In some of these hospitals, as well as in community mental health centres, nurses with graduate degrees share in the treatment of patients and conduct group therapy.

Public health nursing. Public health, or community health, nursing may be either governmental or private. Visiting-nurse services care for the sick at home and carry on individual, family, and community programs of prevention and health teaching. These services are usually supported by civic and other forms of private philanthropy and by fees from patients who can afford to pay. Some of the visiting-nurse services have contracts with local government health and welfare agencies and with industrial health units to provide services. Many have deep, traditional roots in their communities.

Support

In 1978 WHO declared its goal of "health for all" by AD 2000, stressing the efficient and effective use of medical personnel and urging the greater use of nurses in primary care. WHO called for the reorientation of the basic nursing curriculum from an institutional to a community basis, concentrating on education and prevention. Many member states have made progress in that direction. Nursing schools in such countries as Chile, Colombia, Costa Rica, Cuba, Ecuador, Honduras, Indonesia, Nicaragua, Panama, Peru, the Philippines, and Thailand have increased their emphasis on the theory and practice of community health. Other member states, including several in Africa, have begun training for an expansion of nursing functions to strengthen primary health care services, producing nurse practitioners licensed for diagnostic and therapeutic activities, previously considered the domain of the physician, and the prescription of drugs. These functions are usually performed in health centres and in health posts within an organized administrative framework. In some countries nurses are often the only health professionals in isolated rural areas, and even in urban areas specially trained nurses can save the physicians' time and allow for the provision of better service.

Nurses as educators. Nursing education is a field that combines nursing with the teaching of students of nursing and, for some, with the administration of educational programs. A high proportion of the teachers in nursing-education programs teach in clinical situations, in which students learn to care for patients and families in hospitals, at home, and in other situations. Teaching by nurses in staff-development programs of hospitals and other health agencies and in continuing-education programs is usually offered by universities or associations.

Private, office, and industrial nursing. Private, office, and industrial (occupational health) nursing as fields of employment account for only a small percentage of registered nurses. Licensed practical nurses are also employed in these fields.

Military nursing. Military nursing provides an essential part of the health care given to men and women of the

armed services in most countries. During World War I many of the trained nurses who were part of the military forces of the countries involved were drawn from a reserve maintained by the Red Cross. By the close of World War II much progress had been made in assigning nurses rank and responsibilities commensurate with their training and abilities. Medical corpsmen today save lives on battlefields and hospitals and are highly skilled. Some remain in health work after they leave military service.

Government nursing. This field includes military nursing and public health nursing as carried on by national and local health departments, functions discussed above.

Education. The basic educational program for nurses in many countries is both scientific and humanistic in content. In the United States most programs lead to the bachelor's degree. Nursing specialists, teachers, and other leaders in the field, however, may need advanced training at the master's or doctoral level. All educational programs include experience with patients in hospitals, homes, or other settings.

In almost all countries with nursing education there are at least two kinds of programs—those leading to diplomas and those that train auxiliaries, though a large portion of auxiliaries in some countries are untrained. A growing number of countries have one or more bachelor's degree programs; some have several. Among the latter are Australia, Canada, Colombia, Peru, Egypt, India, the Philippines, Taiwan, Thailand, the United Kingdom, and the United States. Master's degree programs have also been developed in several countries.

Countries
with
bachelor's
degree
programs

Post-basic programs for nurses with diplomas have been established in many countries. Some programs offer courses in general education, as well as nursing courses, and some, in universities, may become programs leading to a bachelor's degree. The purposes of such programs vary and include the preparation of teachers, supervisors, or administrators and of nurse specialists in various fields, including midwifery, public health, and teaching of auxiliaries. Some augment the education received in other programs. Enrollment is generally small in relation to the need for their graduates.

The development of nursing education in any country is affected by the developments in general education. In the United States and some other countries, for example, high school graduation or its equivalent has for many years been a requirement for admission to schools preparing registered nurses. In the United States this is also a requirement for admission to practical-nurse programs. In some countries fewer years of previous education are required.

Licensing and registration. After the number of nurses has become substantial and the essential nature of nursing has become established in a country, the need to regulate the practice of nursing under law becomes evident. These laws are aimed at the protection of the public.

Laws define nursing and establish titles under which nurses practice. Most laws establish boards empowered to give examinations and maintain registers of qualified nurses. Boards also are usually empowered to maintain lists of schools the graduates of which are eligible to take examinations. Further, the laws that regulate nursing usually provide ways to recognize licenses acquired by nurses from other states and countries.

The Nurses Act of 1919 established the General Nursing Council for England and Wales to maintain a register of nurses; similar acts were passed in 1919 for Scotland and in 1922 for Northern Ireland. The advent of the National Health Service (1946) called for enlargement of the council and its functions, and in 1949 an act was passed that consolidated all previous acts and established regional committees to work with schools of nursing for their improvement and for ways to meet new challenges.

In 1979 the Nurses, Midwifery and Health Visitors Act set up the Central Council for Nursing, Midwifery and Health Visiting and national boards for the four parts of the United Kingdom. The Central Council became responsible for the basic course, the registration of nurses, disciplinary machinery, and all post-basic education, including health visiting, district nursing, and school nursing.

In the United States, national and state nursing associa-

tions are prime movers in urging legislation for licensure and registration and in securing amendments. The American Nurses' Association gives leadership to state associations and state boards regarding changes in definitions, the desirable provisions of the law, and the legislative process.

The International Council of Nurses exerts leadership among national nursing associations in their efforts to pass and amend laws. The staff of the World Health Organization also has assisted new countries in the licensing process.

Organizations. Near the beginning of the 20th century, nurses began to organize national associations. The purposes of such organizations usually include the promotion of nursing care of high quality, promotion of desirable legislation in nursing and health, formulation of nursing and educational standards, professional development and welfare of nurses, and representation of nursing with other professional associations and government agencies. Many nursing associations publish professional journals. In large countries the national associations have state or provincial and district constituent associations.

An example of a thriving and vigorous national organization is the American Nurses' Association, founded in 1896, which operates aggressive programs. It affects legislation related to health and nursing in the U.S. Congress through its educational and lobbying activities and exerts leadership in the revision of state laws governing registration and the practice of nursing. A strong program for economic security has elevated nurses' salaries and the conditions of patient care, partly through collective bargaining. The official statements of the association set the course of action. Its official publication is the *American Journal of Nursing*. In England a large national nursing organization is the Royal College of Nursing.

In 1899 the International Council of Nurses was founded as a federation of autonomous national nursing associations. Today there are about 100 member associations from most major countries, and the council is working in other countries to assist in the development of national associations and ultimate membership in the council and to encourage national legislation on registration and nursing practice. An international congress is held every four years. The council's journal, *International Nursing Review*, makes a substantial contribution to worldwide literature of nursing. The Florence Nightingale International Foundation, which is related to the council, has conducted an international conference on research in nursing, as well as an international exchange of nurse scholars.

Several national colleges of nurse-midwifery and an international college work for the improvement of maternal and infant health. Nurses participate in international and national associations for public health, mental health, industrial health, and school health and in such areas as heart-disease, cancer, and respiratory-disease control. The Nursing Section of the American Public Health Association is one of its largest sections. Practical nurses and auxiliaries are usually organized, if at all, in separate organizations. Various groups form special organizations in several countries—deans of schools of nursing and surgical nurses are examples. Nursing councils are parts of regional commissions on higher education in the United States.

A unique organization—the National League for Nursing—combines nurses, related professionals, and other public-spirited citizens in the United States for meeting the nursing needs of people throughout the country. This organization also carries on a program of national accreditation of the various kinds of educational programs. Its journal is *Nursing & Health Care*.

ROLES OF INTERNATIONAL AND MULTINATIONAL ORGANIZATIONS

The International Red Cross. The International Red Cross plays two major roles in nursing: it affords educational opportunities for nurses, and it affords nurses opportunity to serve in programs embodying the Red Cross principles of humanity, impartiality, and neutrality. Through participation in Red Cross activities, public-spirited citizens acquire valuable knowledge of the health needs of people and how nursing helps to meet these

needs. Also, nurses broaden their concepts of community and worldwide service.

Some national Red Cross societies operate schools of nursing, post-basic educational programs, and training programs for auxiliaries. In addition, many societies carry on programs preparing nurses to teach health practices in the home; this instruction reaches many thousands of people, who thus learn individual responsibility for health and the care of mothers and children and of patients suffering from simple illnesses. Many societies also train nurses to care for people in disasters and emergencies and deploy nurses to sites of emergencies when needed; some nurses volunteer to serve in emergencies that occur in other countries.

National Red Cross societies place emphasis on community development and youth programs. Societies are urged to be constantly alert to social changes, such as those resulting from urbanization, migration, or drug abuse, and to coordinate their efforts with other voluntary and official agencies. Nurses with knowledge of community conditions can and do contribute to planning and action through board and community memberships, as well as through their participation in the carrying out of Red Cross programs.

The World Health Organization. The World Health Organization has included nursing in its activities from its beginning in 1948. Nurses are included on teams such as those concerned with maternal and child health, malaria, and tuberculosis. Member nations request assistance in developing educational programs for nurses, auxiliaries, and midwives and to organize public health programs and hospitals.

Countries are assisted in establishing nursing as a part of their national health departments to assure the planning and implementing of the nursing portions of programs. Governments are aided in the establishment of nursing and nursing-education systems, including those in midwifery. Other programs in which assistance is given include upgrading diploma schools; development of basic and post-basic programs in universities; revision of entry requirements; coordination of classroom teaching with the clinical practice of students; adaptation of hospitals and health agencies for students' experience; preparation of public health nurses, administrators, midwives, and teachers, including teachers of auxiliaries and of indigenous midwives. (In many countries more than half of all the births are attended by untrained midwives.)

Fellowships are granted to nurses for overseas study, mostly in the areas of teaching, administration, public health nursing, midwifery, maternal and child health, and for learning to plan health services. When consultants from the World Health Organization work in countries, they strive to leave national counterpart personnel to continue the work. Study outside the country may be needed to develop such personnel.

The European Economic Community. In countries of the European Economic Community, nursing is bound by the directives and regulations published by the Council of Ministers in 1977. The Permanent Committee of Nurses lays down broad guidelines for nurses' training and education and each country has a central authority, such as the Central Council of the United Kingdom. In 1979 the nursing directives required universal recognition of certificates and diplomas, and the European Economic Community countries now accept the competence of nurses trained in member states. Also established were a standard of general education to be attained before entry to a school of nursing, the number of hours the training must occupy, and the areas in which students must have theoretical and clinical knowledge. (L.P.Lc./Ed.)

WHO
assistance

The Inter-
national
Council of
Nurses

Dentistry

Dentistry is that profession concerned with the prevention and treatment of oral disease, particularly disease of the teeth and supporting tissues. In addition to general practice, dentistry includes many specialties and subspecialties, including oral surgery, prosthodontics, periodontics, orthodontics, pedodontics, and public health.

HISTORY OF DENTISTRY

Accounts of toothache and other oral complaints together with suggested forms of treatment have been found in some of the earliest texts of the ancient civilizations of the Middle and Far East, including China, Egypt, India, and Mesopotamia. Sumerian clay tablets probably of 5000–3000 BC record the belief, which persisted for several millennia, that dental decay and toothache were caused by the gnawing away of the tooth by a minute worm. In India a professional class of physicians who included dentistry in their activities described teething as a cause of serious infantile illnesses such as cough, diarrhea, vomiting, and convulsions, a myth that has persisted among some people to the present day.

Ancient remedies

In these early civilizations, a great variety of dental medications and operations, such as filling or extracting decayed teeth and the splinting of loose teeth or fractures of the jaws, were in use. In about 400 BC Hippocrates described many oral diseases, and he is credited by many with having introduced the term *apthae*, which is still in use for the painful, but otherwise harmless, oral ulcers that occur particularly in adolescence or early adult life. Another observation made in ancient Greece, that sweet foods caused teeth to decay, was only experimentally confirmed in the 20th century.

Over the centuries many kinds of dental instruments and empirical forms of treatment developed, but dentistry remained essentially a craft included in the activities of surgeons or practiced by itinerant “healers.” In the early 16th century in France, dentistry was practiced by barbers, and in England the Barber-Surgeons’ Company (later to become the Royal College of Surgeons) was granted a charter by Henry VIII. Surgery and dentistry eventually passed out of the hands of barbers, and the 16th and 17th centuries saw the emergence of dentistry as a specialty with its own literature.

The first textbook in dentistry was published in Leipzig in 1530, and 50 years later students of dentistry were admitted to the University of France. By 1622 a number of men had been granted the title of surgeon-dentist, although the title was not fully established for a number of years after that. During the reign of Louis XIV, the surgeon-dentists formed a separate subdivision of the surgeons’ guild and, the year after the subdivision was formed, it became law that those who wished to practice in the field of mouth surgery and artificial restoration had to pass prescribed examinations. At this time some women were permitted to practice dentistry in France, although the privilege was revoked during the mid-1700s.

In the 17th century in England dentistry was referred to as an independent specialty. The first English text on the subject of dentistry was published in 1685.

Textbooks

Publication of three highly significant textbooks occurred during the 18th century. The French dentist Pierre Fauchard’s *Le Chirurgien Dentiste* (“The Surgeon Dentist”) in 1728 put dental treatment on a more scientific plane than ever before and advocated a broader education for dentists. In 1756 the German dentist Philipp Pfaff’s *Abhandlung von den Zähnen* (“Treatise on the Teeth”) appeared, and in 1771 an anatomist and surgeon of England, John Hunter, who was giving lectures on dentistry, published *The Natural History of the Human Teeth*. Joseph Fox was appointed dental surgeon at Guy’s Hospital in 1799, and at the same time lectures on dentistry were set up in Guy’s and in a number of other areas.

The next major advance in dentistry was the introduction of general anesthesia to medicine by two dentists, Horace Wells and William Morton, between 1844 and 1846. In 1884 Sigmund Freud demonstrated the value of cocaine as a local anesthetic. William Stewart Halsted, a surgeon, carried out the first regional local anesthetic of the inferior dental nerve, using a technique that has remained essentially unchanged.

In 1890 W.D. Miller published his finding that teeth could be caused to decay, *in vitro*, by the action of mouth bacteria on carbohydrates. More than half a century was to elapse, however, before the bacterial basis for dental caries could be definitively established. Since the 1940s, dental research has increased in range and scale in all of

its specialties. The result has been an increasing number of publications and journals and the formation of national and international research associations.

No formal dental schools were established in England until 1858, although dental hospitals had been established earlier. These dental hospitals, which were created as centres to serve the poor, were founded and supported by dentists. For the most part they were independent of medical schools and general hospitals. They did accept some students, however, who partly provided cheap labour to operate the hospital and whose fees were used to support the charitable work undertaken. The honorary dental surgeons were responsible for any teaching that was done. In 1858 the first dental school in the United Kingdom was established by the Odontological Society of London and the second school the following year by the College of Dentists of England. Both were private schools.

Early dental schools

At about the same time the Royal College of Surgeons arranged to hold examinations for licensure for dental surgery. It remained the examining body for dentistry in Great Britain for about 20 years. In 1878 the first Dentists’ Act in the United Kingdom was passed and the General Medical Council established a register of those qualified to hold the title of dental surgeon. The act did not, however, prohibit individuals without these qualifications from practicing. The General Medical Council also prescribed a curriculum for the training of dentists. It required two years of preceptorship to learn dental mechanics and three years in a medical school and dental hospital. An inspector of the examination procedures at the colleges in 1897 recommended that the one examination be replaced by three—in preliminary science and dental mechanics and a final examination. These recommendations were put into effect and remained without much alteration until 1922.

While schools were being established and standards for licensure created in England, development had also been going on in North America. Dentistry was being practiced in the United States by the late 1790s, and the first textbooks appeared in the early 1800s. By 1834, when the first meeting of a dental society in the United States was held, services were being provided by three distinct classes of people. First, there were those who had qualified by a course of study in the principles of medicine and surgery. Second, there were those who had obtained a preparatory course of medical study and then begun practicing dentistry without having studied the mechanics of dental operations. These people were held in comparative respect since in due course they did obtain some degree of skill in their operations, although they were criticized for not having obtained this skill before actually inflicting their services on the public. The third group included a great number of charlatans—anyone who decided he would like to practice dentistry. It is this group that brought dentistry into disrepute.

Early dentistry in the United States

During the 1830s two unsuccessful attempts to establish dental schools were made, the first in Kentucky and the second in New York City. The latter failed to get under way because of staff problems. No money was available at that time for salaries, and local dentists evidently felt that they could not spare time from their practice for teaching.

The first dental school in North America was established in Baltimore in 1840. Subsequently many other schools were started, many of which were operated privately or commercially. All of the early schools were separate from universities. The first, and for many years the only, school of dentistry in the United States associated with a university was opened at Harvard in 1867.

This was also the era during which dentistry became organized in Canada. The first dental act in Canada was passed by the legislature of the Province of Ontario in 1868. This act incorporated the Royal College of Dental Surgeons of Ontario and gave it the dual function of teaching and licensing. License requirements stipulated five years of practice in a dental office for registration. Except for a simple act that was passed in Alabama in 1841, the Ontario Act was the first law respecting dental practice in North America. The regulations for registration were fairly simple by today’s standards, but with the large number of itinerant dentists, who had no formal training

whatever, moving back and forth across the Canadian-U.S. border at that time, even this five-year requirement proved useful.

THE PRACTICE OF DENTISTRY

Licensure requirements. The practice of dentistry is now well controlled, and in all countries of the world in which dentistry is practiced there is a licensing requirement. The licensing authority may be the government or national dental organizations.

Canada. In Canada each province has its own licensing authority. This can be a college, such as the Royal College of Dental Surgeons of Ontario, or an association, such as the Manitoba Dental Association. There is also a national authority, the Dental Examining Board.

The university degree (doctor of dental surgery or doctor of dental medicine) does not in itself entitle the holder to practice but is an academic qualification for presentation to the licensing board under whose jurisdiction the holder wants to obtain a license to practice. The regulations of the provincial licensing boards vary but usually require an examination for licensing.

United States. Licensing authority in the United States is vested in state boards of dental examiners, most of which require an examination. Most states require U.S. citizenship as a prerequisite. Some states require non-citizens to submit the declaration of intent to become a citizen or the first papers as a requirement for admission to the dental licensure examination; a few do not require citizenship for such admission.

Nationals with foreign diplomas may be admitted to practice if their diplomas were issued by a school approved by the American Dental Association and if they pass the state licensure examination.

Soviet Union. Dentists in the former Soviet Union were divided into three distinct classes, known as two-, three- and five-year dentists. The two-year dentists were dental technicians who studied for two years beyond secondary school, after which they became eligible to work in a dental laboratory. Two-year dentists were not permitted to treat patients. The three-year dentists were permitted to treat patients, but their practice was limited to restorative, prosthetics, and prophylactic dentistry. The five-year dentists were stomatologists, who were on a par with physicians in the Soviet Union. They received training equal to that of physicians and, in addition, were trained in operative dentistry, crown and bridge dentistry, prosthetics, exodontia, and anesthesia, both general and local. The five-year dentist's degree entitled the practitioner to perform surgery on the hard as well as all soft tissues in the mouth and throat.

European Economic Community. With the advent of the European Economic Community and the Council of Europe, it has become accepted policy that doctors and dentists should be able to move freely and practice within any of the member countries. For this to be acceptable there has had to be mutual recognition of dental degrees and comparable forms of qualification. The Council of the European Communities has, therefore, agreed on a set of directives that set out the training requirements for dental education in the member states. This has created no difficulties for most European countries, where dentistry has long been recognized as a specialty in its own right. In Italy and Spain, however, where dentistry is a subspecialty of medicine, transitional provisions have had to be made until dental training can be harmonized with that of other member states.

Permission to practice in the United Kingdom is granted by the General Dental Council to those holding (1) a degree or diploma in dentistry or dental surgery conferred in Great Britain or Northern Ireland, (2) a degree or diploma in dentistry or dental surgery granted elsewhere that has been recognized by the General Dental Council, or (3) a degree or diploma approved by the General Dental Council provided that these graduates have passed the statutory examination written under arrangements made by the General Dental Council.

Dentists in Germany must hold a dental surgeon's diploma, which authorizes private practice without fur-

ther examination. They must be registered by local dental boards and by health authorities.

In Italy a diploma in dentistry, which allows the use of the title of Specialist in Diseases of the Mouth, Teeth, and Jaws, constitutes a license to practice. Holders of the diploma of Doctor of Medicine have passed examinations in dentistry and for this reason may also practice dentistry but do not have the specialist title.

Japan. Since about 1903 Japanese dentistry has been mainly patterned after that practiced throughout the United States. Those wanting to practice dentistry or dental surgery must be recognized by the national government. Applicants for registration must pass the national examination for dentists and obtain license to practice. These requirements must also be fulfilled by registered medical practitioners wanting to practice dentistry, by Japanese citizens, and by foreigners who have qualified in Japan.

Types of practice. *Private practice.* In Canada, the United States, the United Kingdom, and Australia, dentists in private practice constitute the vast majority of all licensed dentists. The situation is much the same in France and various other countries.

Dental practice has changed significantly since 1920, without a concurrent change in the basic dental curriculum. Dental procedures have shifted from the repair and extraction of teeth for the relief of pain in 1920 to prevention of disease. Dental practice has also changed in larger urban centres from the isolated private practice common in 1920 to a complex system of groups of professionals in a central location. Extensive use is now made of dental hygienists, who often receive the patient from the examining dentist. Dental hygienists provide services such as preventive procedures, prophylaxis, scaling, X rays, and dental health education. Most practices also use dental assistants.

Another development that has occurred in dental health-care services is the extension of the duties currently carried out by dental auxiliaries. New Zealand has pioneered in the field with the creation of the dental nurse, an auxiliary who is trained to provide dental care for children without the supervision of a dentist. The United Kingdom has also developed the "dental auxiliary," who performs somewhat similar duties but under closer supervision. In Canada and the United States, pilot projects have been conducted to test the feasibility of using dental auxiliaries for certain operative procedures in order to increase productivity, quality, and general service to the public.

France may be taken as an example of the development of the practice of dentistry in continental Europe. There are two types of dentists practicing in France, the *chirurgien dentiste* and the stomatologist. The practice of dentistry in France by a *chirurgien dentiste* has, since 1892, been restricted to persons of French nationality who hold a state diploma and who are registered with the Order of Dentists. The Order of Dentists is responsible for registration and discipline but is not concerned with dental education, which is controlled by the state through the common state diploma.

Stomatologists are practitioners who have a diploma in medicine together with either a diploma in dental surgery or a certificate of special studies (two years) in stomatology. Specialization within the field of dentistry is not encouraged. There are no rules laid down for it nor are there any special courses or diplomas or titles.

Hospital dental practice. Three types of dental care are normally carried out in the hospital environment: (1) clinical procedures normally provided in a dental office, for ambulatory inpatients and outpatients, (2) bedside care for persons admitted for other medical reasons, and (3) inpatient care for patients admitted to a hospital for purely dental conditions.

Dentists may treat patients in hospitals either privately, on a fee-for-service basis, or under some form of government program, such as the National Health Service in the United Kingdom or the Provincial Medicare Plan (surgery only) in Canada. Hospital dental services have for years been an integral part of dental health care and dental education in the United Kingdom, and such services by hospital dental departments have expanded steadily in the United States and Canada.

Shift from individual to group practice

Types of dental care

Hospital dental departments are normally established in the same manner as any other hospital department and are headed by a chief of service, who has the same status as other chiefs of service within the hospital. In some instances, the chief of the dental department may be responsible to the chief of surgery. There are two types of hospital dental departments—one that is established in a teaching hospital and the other in a general hospital with no teaching component. In the teaching hospital the dental department is associated with a faculty of dentistry and forms an integral part of the undergraduate curriculum and, if they exist, of the graduate and postgraduate programs. One of the chief purposes of hospital dental departments is to make available the service of consultants to other hospital departments and general practitioners. This service is most highly developed in teaching hospitals. Usually, certain general dental treatment is provided for inpatients and outpatients. Hospital dental services or departments are prevalent in western Europe.

Public health practice. Generally typical of dental public health practice in Canada and in many areas of the United States is the program carried on in Ontario. There dentists trained in public health, hygienists, and dental assistants carry out a preventive and educational program basically concerned with the examination of children, the recording of basic dental conditions, and the provision of dental health education.

Military practice. Most countries of the world provide dental-care service for their armed forces. The organization of such a service varies extensively. In Canada the Royal Canadian Dental Corps has the same status as the Royal Canadian Medical Corps, with a brigadier general as the director. Military service for dentists in the United States is under the U.S. Public Health Service, the chief of service being an assistant surgeon general. In the United Kingdom dental care is provided by three separate dental branches—Navy, Army, and Air Force.

Governmental practice. In many countries dentists are required to work a number of years for the government before they may be considered private practitioners of the type known in Canada and the United States. This service requirement may be based on the fulfillment of an obligation for government financial support during undergraduate training, or there may be a government regulation that all dental graduates must work for the state for a prescribed number of years. Another example of government practice is in the United Kingdom, where dentists are employed by local authorities to provide dental care under the Maternal and Child Welfare Services and the School Dental Service.

The employment of dentists on a salary basis for the general practice of dentistry is not extensive in the United States or Canada. At the national level it may be the provision of dental care for eligible Indians and Eskimos, war veterans, or inmates of penitentiaries. At a municipal level, dentists may be employed in a school dental service. Dentists in both Canada and the United States commonly agree to provide service for families who qualify for social assistance. They are paid on a fee-for-service basis; the fee schedule is usually set, normally after consultation with the profession, by the agency responsible for the social service plan.

Government medical care was introduced in Japan in the late 1930s. This system was expanded until by 1962 almost the entire population was covered. There are limitations to the services offered by government medical care, as in orthodontics or in preventive dentistry.

Dental specialties and subspecialties. In most countries that recognize specialties in dentistry, the specialist is limited to practice in the specialty and cannot carry out the practice of general dentistry. Where the specialty is thus limited the general dentist may refer patients, and a specialist's practice is mainly on a referral basis. In Britain and in certain provinces in Canada, specialists may conduct a general practice.

Orthodontics. Orthodontics takes as its aims the prevention and correction of malocclusion of the teeth and associated dentofacial incongruities. Orthodontics has been practiced since the days of the ancient Egyptians,

but methods of treatment involving the use of bands and removable appliances have become prominent only since the beginning of the 20th century. The United States gave impetus to the development of orthodontics, which was recognized as a specialty with the formation of the American Society of Orthodontists in 1900.

The demand for this service extends from the child to the mature adult, although human bone responds to tooth movement best in a person under 18, and it is generally agreed that children benefit more from treatment than do adults. In general, oral health and physical appearance are the two most important reasons for undertaking a course of orthodontic care.

Pedodontics. Pedodontics, analogous to pediatrics in medicine, is concerned with the dental care of children and adolescents.

Much of the routine of practice is centred around the control of caries (tooth decay) and involves the use of fluoride and dietary and hygienic instruction. The need to influence tooth positions presents the next most frequently encountered problem. The correction of incipient abnormalities in tooth alignment may obviate the necessity for lengthy treatment. Many pedodontists use growth-influencing techniques to correct jaw alignments. A working knowledge of children's behaviour patterns, patience, and a knowledge of childhood physical and mental diseases and their ramifications are important qualifications of the pedodontist.

Periodontics. Periodontics is concerned with the prevention, diagnosis, and treatment of diseases of the periodontal tissues—the tissues that surround and support the teeth. These tissues consist mainly of the gums and the jaw bones and their related contiguous structures.

The most prevalent periodontal disease is periodontitis, commonly called pyorrhea, an inflammatory condition usually produced by local irritants. Periodontitis, if untreated, destroys the periodontal tissues and is a major cause of the loss of teeth in adults.

The advances of periodontics have been mostly in techniques of treatment. It is believed that bacterial plaque, a soft layer of substances rich in bacteria that adheres to the teeth, is the factor responsible for most destruction of the gums and the tissues surrounding the teeth. Periodontists advocate removal of such plaque by a specific regimen of controlled hygiene.

Prostodontics. Prostodontics is concerned with the restoration and maintenance of oral function, comfort, appearance, and health by the replacement of missing teeth and contiguous tissues with artificial substitutes, or prostheses.

There are three main branches of the specialty, concerned, respectively, with removable prostheses, fixed prostheses, and maxillofacial prostheses. Maxillofacial prostheses are supplied to persons who have suffered congenital, traumatic, or surgical defects of the mouth, jaws, or associated facial structures.

The proper fitting of oral prostheses requires a detailed knowledge of the anatomy of the head and neck, of the physiology of the neuromuscular system, and of the science of occlusion and jaw movements. It also requires skill in planning, mouth preparation, impression making, registration of jaw relations, try-in procedures, placement of the prostheses, and follow-up care.

Oral medicine. Oral medicine, or stomatology, treats the variety of diseases that affect both the skin and the oral mucous membranes. Some of these diseases, such as pemphigus vulgaris, can develop their first manifestations in the mouth and can be life-threatening. Cancer of the mouth also has a high mortality rate, partly because it grows in such close proximity to so many vital structures and readily involves them. With all such diseases of the oral cavity, removal of a portion of the lesion for examination under the microscope (biopsy) by an oral pathologist is an essential procedure, and many other laboratory procedures are often also required for the diagnosis of oral mucosal diseases.

Oral pathology. Oral pathology is the study of the causes, processes, and effects of oral disease, together with the resultant alterations of oral structure and functions.

Dental care for children

Branches of prostodontics

The status of the specialist

The oral pathologist provides diagnoses on which treatment by other specialists will depend.

Oral surgery. Oral surgery deals with the diagnosis of, and the surgery required by, the diseases, injuries, and defects of the human jaws and associated structures. Both dentists and physicians refer a wide variety of special dental problems to the oral surgeon. These may include the removal of impacted and infected teeth and the treatment of cysts, tumours, lesions, and infections of the mouth and jaws. In addition there are more complex problems, such as jaw and facial injuries, cleft palate, and cleft lip.

Public health dentistry. Dental public health is recognized as a specialty in Canada and the United States. The American Dental Association recognizes dental public health as a specialty if the holder of the master's degree proceeds to a further year of study in training and passes the examination of the American Board of Dental Public Health. Training in dental public health is available in the United Kingdom, but the specialty is not emphasized to the same degree in the rest of the world.

Forensic dentistry. Forensic dentistry is the study and practice of aspects of dentistry that are relevant to legal problems. It is a specialty practiced by few and is not usually part of dental education. Forensic dentistry is, however, of considerable legal importance for several reasons, one of the most important of which is the fact that the teeth are the structures of the body most resistant to fire or putrefaction. Moreover, the arrangement of the teeth or any restoration in them is virtually or completely unique to any given individual and, if dental records can be found, may enable identification with certainty similar to that provided by fingerprinting. The identification of human remains after aircraft accidents, for example, can often be made only by this means. Minor irregularities of the teeth can also be reproduced in bite marks, enabling a suspect to be identified if he has bitten another person.

Dental education. *Predental programs.* In a majority of countries in the world undergraduate training in dentistry is available. Many require predental training prior to acceptance into a school of dentistry. The predental training is in addition to primary and secondary education, which usually takes from 10 to 12 years. The required number of years in predental education varies from one to seven (a number of European countries require from five to seven years of medical education before entering dentistry). Predental course training usually includes such studies as biology, chemistry, physics, and mathematics. Certain faculties of dentistry in Canada and the United States require a bachelor's degree in arts or science as a prerequisite for admission into a dental faculty.

Dental school and training. After predental courses, training consists of four years in a faculty of dentistry to qualify as a doctor of dental surgery (D.D.S.) or doctor of dental medicine (D.M.D.). The program of studies during the four-year course includes the following biological sciences: human anatomy, biochemistry, bacteriology, histology, pathology, pharmacology, microbiology, and physiology, upon which the succeeding studies of the theory and techniques of dental practice are based. Studies required with respect to dental practice include restorative dentistry, prosthetics, orthodontics, surgery, preventive dentistry, medicine, dental public health, pedodontics, periodontics, radiology, clinical practice, and anesthesia.

ANCILLARY DENTAL FIELDS

Dental hygienists. The hygienist is a figure in the campaign to reduce periodontal disease and to improve physical well-being by promoting better care of the mouth.

The prevention of oral disease through education and treatment is the chief function of hygienists. The specific duties and services that they are allowed to perform depend on the bylaws of the licensing bodies, the requirements of the dental offices in which they are employed, or the aims and objectives of the public health programs in which they are engaged. At all times hygienists work under the effective supervision of a qualified dentist, and they are not permitted to establish their own practice.

Hygienists employed in dental offices remove deposits and stains from the patient's teeth, apply fluorides, and

observe and record conditions of decay and disease for the dentist's information. Further duties may include the taking of X-ray photographs of parts of the mouth, which the hygienist develops and mounts. Another function of the hygienist is to promote dental health by advising on diet and nutrition and encouraging oral hygiene.

Hygienists employed by educational authorities assist school dentists by performing such duties as examining children's teeth. They may also visit classrooms to explain the importance of oral hygiene and to give instruction in the proper care of the teeth and gums. In hospitals they perform mainly the same duties as for private practitioners.

Dental nurses and dental auxiliaries. In New Zealand an auxiliary known as the dental nurse has been carrying out a dental-care program for children for a number of years. The dental nurse receives minimal supervision but is equipped to provide a dental-care program for children up to 13 years of age. The dental nurse is trained for two years in a special course for dental nurses with entrance requirements below the university level.

Dental assistants. The majority of dentists in private practice employ one or more dental assistants to provide such services as the reception of patients, the keeping of records and accounts, assistance for the dentist while he is treating patients, general upkeep of the office, developing of dental X rays, and sterilization of instruments.

Dental technicians and dental mechanics. Dental technicians, also called dental mechanics, make artificial crowns, bridges, dentures, and other dental appliances according to dentists' specifications. Work orders, accompanied by models or impressions of patients' mouths, state the exact requirements for each particular job. In large laboratories the various stages of manufacture are often divided and the technicians employed may specialize. Sometimes partially skilled persons are hired to work in limited aspects of production on an assembly-line basis.

ORGANIZATIONS

Associations of dentists, dental journals, and dental schools now exist in almost every country of the world. The Fédération Dentaire Internationale (International Dental Federation) was founded in 1900 and has met annually except in times of war. It has sponsored international dental congresses that are planned to meet every five years. Other international organizations include the International Association for Dental Research (Association Internationale pour la Recherche Dentaire) and the Association pour les Recherches sur les Parodontopathies (Association for Research into Periodontal Diseases), which was organized in 1932. *The International Dental Journal*, published by the Fédération Dentaire Internationale, was founded in 1950.

Within the general framework of the World Health Organization, the dental health program has progressed steadily from the beginning. A proposal for a joint review of stomatology and dental hygiene in collaboration with the International Dental Federation was made at the first World Health Assembly in 1948.

Certain organizations, including the World Health Organization and Fédération Dentaire Internationale, and countries such as New Zealand and the United States offer direct and financial assistance to many developing countries in the development of health educational and dental care services. For example, New Zealand has long provided developing countries with the benefit of its experience in the use of dental auxiliaries or what is commonly known as school dental nurses. Direct assistance is provided in the development of other public health dental services to countries such as Sri Lanka, Malaysia, Singapore, Brunei, Thailand, Indonesia, Hong Kong, and Papua New Guinea. Dentists from these countries have had the opportunity of studying the New Zealand system and a number of school dental nurses have received their training there, enabling them to assist in the establishment of their own training facilities. (R.A.Co.)

Pharmacy

Pharmacy is the science and art concerned with the preparation and standardization of drugs. Its scope includes the

Training and functions of dental nurses

Chief function of hygienists

cultivation of plants that are used as drugs, the synthesis of chemical compounds of medicinal value, and the analysis of medicinal agents. Pharmacists are responsible for the preparation of the dosage forms of drugs, such as tablets, capsules, and sterile solutions for injection. They compound physicians', dentists', and veterinarians' prescriptions for drugs. The science that embraces knowledge of drugs with special reference to the mechanism of their action in the treatment of disease is pharmacology.

HISTORY OF PHARMACY

The beginnings of pharmacy are ancient. When the first person expressed juice from a succulent leaf to apply to a wound, this art was being practiced. In the Greek legend, Asclepius, the god of the healing art, delegated to Hygieia the duty of compounding his remedies. She was his apothecary or pharmacist. The physician-priests of Egypt were divided into two classes: those who visited the sick and those who remained in the temple and prepared remedies for the patients.

In ancient Greece and Rome and during the Middle Ages in Europe, the art of healing recognized a separation between the duties of the physician and those of the herbalist, who supplied the physician with the raw materials from which to make medicines. The Arabian influence in Europe during the 8th century AD, however, brought about the practice of separate duties for the pharmacist and physician. The trend toward specialization was later reinforced by a law enacted by the city council of Bruges in 1683, forbidding physicians to prepare medications for their patients. In America, Benjamin Franklin took a pivotal step in keeping the two professions separate when he appointed an apothecary to the Pennsylvania Hospital.

The development of the pharmaceutical industry since World War II led to the discovery and use of new and effective drug substances. It also changed the role of the pharmacist. The scope for extemporaneous compounding of medicines was much diminished and with it the need for the manipulative skills that were previously applied by the pharmacist to the preparation of bougies, cachets, pills, plasters, and potions. The pharmacist continues, however, to fulfill the prescriber's intentions by providing advice and information; by formulating, storing, and providing correct dosage forms; and by assuring the efficacy and quality of the dispensed or supplied medicinal product.

THE PRACTICE OF PHARMACY

Education. The history of pharmaceutical education has closely followed that of medical education. As the training of the physician underwent changes from the apprenticeship system to formal educational courses, so did the training of the pharmacist. The first college of pharmacy was founded in the United States in 1821 and is now known as the Philadelphia College of Pharmacy and Science. Other institutes and colleges were established soon after in the United States, Great Britain, and continental Europe. Colleges of pharmacy as independent organizations or as schools of universities now operate in most developed countries of the world.

The course of instruction leading to a bachelor of science in pharmacy extends at least five years. The first and frequently the second year of training, embracing general education subjects, are often provided by a school of arts and sciences. Many institutions also offer graduate courses in pharmacy and cognate sciences leading to the degrees of master of science and doctor of philosophy in pharmacy, pharmacology, or related disciplines. These advanced courses are intended especially for those who are preparing for careers in research, manufacturing, or teaching in the field of pharmacy.

Since the treatment of the sick with drugs encompasses a wide field of knowledge in the biological and physical sciences, an understanding of these sciences is necessary for adequate pharmaceutical training. The basic five-year curriculum in the colleges of pharmacy of the United States, for example, embraces physics, chemistry, biology, bacteriology, physiology, pharmacology, and many other specialized courses. As the pharmacist is engaged in a business as well as a profession, special training is provided

in merchandising, accounting, computer techniques, and pharmaceutical jurisprudence.

Licensing and regulation. To practice pharmacy in those countries in which a license is required, an applicant must be qualified by graduation from a recognized college of pharmacy, meet specific requirements for experience, and pass an examination conducted by a board of pharmacy appointed by the government.

Pharmacy laws generally include the regulations for the practice of pharmacy, the sale of poisons, the dispensing of narcotics, and the labeling and sale of dangerous drugs. The pharmacist sells and dispenses drugs within the provisions of the food and drug laws of the country in which he practices. These laws recognize the national pharmacopoeia (which defines products used in medicine, their purity, dosages, and other pertinent data) as the standard for drugs. The World Health Organization of the United Nations began publishing the *Pharmacopoeia Internationalis* in the early 1950s. Its purpose is to standardize drugs internationally and to supply standards, strengths, and nomenclature for those countries that have no national pharmacopoeia.

Research. Pharmaceutical research, in schools of pharmacy and in the laboratories of the pharmaceutical manufacturing houses, embraces the organic chemical synthesis of new chemical agents for use as drugs and is also concerned with the isolation and purification of plant constituents that might be useful as drugs. Research in pharmacy also includes formulation of dosage forms of medicaments and study of their stability, methods of assay, and standardization.

Another facet of pharmaceutical research that has attracted wide medical attention is the "availability" to the body (bioavailability) of various dosage forms of drugs. Exact methods of determining levels of drugs in blood and organs have revealed that slight changes in the mode of manufacture or the incorporation of a small amount of inert ingredient in a tablet may diminish or completely prevent its absorption from the gastrointestinal tract, thus nullifying the action of the drug. Ingenious methods have been devised to test the bioavailability of dosage forms. Although such in vitro, or test-tube, methods are useful and indicative, the ultimate test of bioavailability is the patient's response to the dosage form of the drug.

Licensing systems for new medicinal products in Europe and North America demand extensive and increasingly costly investigation and testing in the laboratory and in clinical trials to establish the efficacy and safety of new products in relation to the claims to be made for their use. Proprietary rights for innovation by the grant of patents and by the registration of trademarks have become increasingly important in the growth of the pharmaceutical industry and its development internationally.

The results of research in pharmacy are usually published in such journals as the *Journal of Pharmacy and Pharmacology* (London), the *Journal of the American Pharmaceutical Association* and the *Journal of Pharmaceutical Sciences* (Washington, D.C.), the *American Journal of Pharmacy* and the *American Journal of Hospital Pharmacy* (Philadelphia), and the *Pharmaceutica Acta Helvetica* (Zürich).

ORGANIZATIONS

There are numerous national and international organizations of pharmacists. The Pharmaceutical Society of Great Britain, established in 1841, is typical of pharmaceutical organizations. In the United States the American Pharmaceutical Association, established in 1852, is a society that embraces all pharmaceutical interests. Among the international societies is the *Fédération Internationale Pharmaceutique*, founded in 1910 and supported by some 50 national societies, for the advancement of the professional and scientific interests of pharmacy on a worldwide basis. The Pan American Pharmaceutical and Biochemical Federation includes the pharmaceutical societies in the various countries in the Western Hemisphere.

There are also other international societies in which history, teaching, and the military aspects of pharmacy are given special emphasis.

(J.C.K./Fr.H.)

Ancient and medieval pharmacy

Bio-availability

The basic five-year curriculum

International organizations

LEGAL ASPECTS OF MEDICINE

Maintenance of professional standards

HISTORY

Physicians historically have set their own standards of care and their conduct has usually been judged by comparing it to that of other physicians. "Ethical" canons or codes generally focused on professional etiquette and courtesy toward fellow physicians rather than on relationships with patients. The Hippocratic oath was a notable exception, but its provisions were only ascribed to by a minority of Greek physicians.

The law has become intimately involved in medical practice only in the 20th century. Until recently legal medicine, or forensic medicine, was a field devoted exclusively to the uses of medicine in the courtroom, primarily in two settings: forensic pathology and forensic psychiatry. The pathologist has traditionally been asked to determine and testify to the cause of death in cases of suspected homicide and to aspects of various injuries involving crimes such as assault and rape. Pathological testimony may also be required in civil cases involving, for example, occupational injury, negligent injury, automobile accidents, and paternity suits. Similarly, when a defendant pleads insanity as a defense, a psychiatrist is asked to examine the defendant and to testify as to his mental state at the time of the crime. The relevant question is usually whether his criminal behaviour was the product of a mental illness or whether the defendant was able to distinguish right from wrong. In civil cases, psychiatrists frequently appear as witnesses in cases of child custody and involuntary commitment.

Since 1960 the legal climate has changed drastically. Civil lawsuits alleging medical malpractice have become a fact of professional life for many Western physicians. Issues formerly relegated to ethics, such as abortion and termination of treatment, also have become important civil rights issues in courtrooms across the world, as have issues of informed consent and patients' rights. Wide-ranging campaigns aimed at arresting the spread of infectious diseases, such as acquired immune deficiency syndrome (AIDS), have involved the legal system in issues of privacy, confidentiality, quarantine, and research using human subjects.

So great has the change been that forensic medicine has now become a subspecialty of a separate field, usually called health law to emphasize its application not only to medicine but to health care in general. This new field of health law is not limited to the courtroom but is active as well in legislatures, regulatory agencies, hospitals, and physicians' offices.

RELATIONSHIP OF LAW AND ETHICS

The legal philosopher Lon Fuller has distinguished between "the morality of aspiration" and "the morality of duty." The former may be denoted ethics, the latter law. Ethics tells people what they should do and embodies the ideals they should strive to attain. Unethical behaviour leads to punishments that are related to how an individual is perceived, both by himself and by his fellow man. Law, on the other hand, provides boundaries of actions, set by society, beyond which a person may go only by risking external sanctions, such as incarceration or loss of a medical license.

This may explain why ethical codes usually involve generalities, while laws tend to be more specific. For example, the Hippocratic oath, formulated in the 5th century BC, is concerned with the physician's doing no harm, refraining from performing abortions and giving deadly drugs, and maintaining strict confidentiality. Law, on the other hand, may permit abortions under certain circumstances, permit the giving of potentially lethal drugs in extreme situations, and sanction the violation of confidentiality when the interests of society demand it.

For example, in treating a patient with AIDS, a physician might administer a potentially fatal experimental drug in a desperate attempt to destroy the retrovirus that causes AIDS. The physician might also warn the patient's spouse

that there is a danger that the virus will be transmitted, even if the patient does not want the spouse to be so informed. These scenarios present ethical dilemmas, and it is unlikely that any of these actions would have legal consequences today.

Although the Hippocratic oath has largely been superseded by such modern oaths as the Declaration of Geneva, the International Code of Medical Ethics, and the Canons of the American Medical Association, these codes of conduct retain the brevity and generality of the Hippocratic oath. For example, the International Code of Medical Ethics, developed and promulgated by the World Medical Association shortly after World War II, provides in part for the following:

A doctor must always maintain the highest standards of professional conduct.

A doctor must practice his profession uninfluenced by motives of profit. . . .

A doctor must always bear in mind the obligation of preserving human life. . . .

A doctor shall preserve absolute secrecy on all he knows about his patient because of the confidence intrusted in him.

Modern advancements in the field of medicine, such as cardiopulmonary resuscitation (which restores regular rhythm to an arrhythmic or failed heart) and mechanical respirators (which breathe for patients unable to expand their lungs) sometimes have been able to postpone a death that previously had been imminent. Under these circumstances, it may be difficult to relate the rules of ethics to the realities of the situation. For example, the meaning of "obligation of preserving human life" becomes unclear in the context of a young woman in a permanent coma who will probably die if the mechanical respirator is removed, but who may live for decades (in a coma) if the machine remains in place. It is not clear that the Hippocratic ideal of doing "no harm" requires that the machine remain in place or that it be removed.

In 1976 these same questions were confronted by the New Jersey Supreme Court in the landmark case of Karen Ann Quinlan. Her parents requested her physicians to remove the mechanical respirator in order to let their daughter die a natural death. The doctors refused, relying primarily on medical ethics, which they believed prohibited taking an action that might lead to the death of the patient.

In court, however, the lawyers for the Quinlan family argued that what was at stake was not medical ethics but the legal rights of the individual patient to refuse medical treatment that was highly invasive and offered no chance for a cure. The court agreed that patients have the legal right to refuse medical treatment, determined that honouring such a refusal was consistent with medical ethics, and decided that the parents of Karen Ann Quinlan could exercise her right to refuse treatment on her behalf. In order to reassure the physicians involved, however, the court also decreed that if a hospital's ethics committee agreed with the prognosis of permanent coma, the removal of the respirator could take place and all parties involved would have legal immunity from civil or criminal prosecution. Karen Ann Quinlan's respirator was removed, although she continued to breathe on her own. She survived in a coma until she died of pneumonia almost 10 years later.

The case of Karen Ann Quinlan has become a parable of modern medicine and of the relationship between medical ethics and the law. Although an issue was made of medical ethics by both the physicians and the court, the case primarily involved medical practice and the fear of potential legal liability. Modern physicians now also worry about the law as well as ethics, and they fear criminal lawsuits that allege homicide or assisting suicide, and civil lawsuits that allege malpractice. To address these concerns, the New Jersey court created an ethics committee with the power to grant legal immunity for actions and to diffuse the responsibility for them.

This model has not been followed by other courts, however, although ethics committees have been established

Effects of modern technology on ethics

Ethics committees

Changes since 1960

The Hippocratic oath

in North America, Europe, and Australia to help educate hospital staffs on such matters as the withholding and withdrawing of treatment and on general ethical conduct with patients. In fact, physicians are rarely taken to court on criminal charges for decisions about patient care that are made in good faith. The Massachusetts Supreme Court, for example, has summarized the criminal law in this regard: "Little need be said about criminal liability: there is precious little precedent, and what there is suggests that the doctor will be protected if he acts on a good faith judgment that is not grievously unreasonable by medical standards."

In the United States, the Quinlan rationale has been expanded to include the right of all mentally competent patients, whether terminally ill or comatose, to refuse any and all medical treatments (including artificial feedings). Some people state their wishes about treatment in documents called "living wills," to which physicians can refer after the patient is no longer able to speak for himself. People also may provide someone with power of attorney to make decisions about medical treatment once they become incompetent to act on their own behalf. The American Medical Association has stated that honoring patients' refusals of treatment is consistent with both medical practice and medical ethics. Other countries, like The Netherlands, have gone further and have held that it is legally and ethically acceptable for physicians to assist a patient with a terminal illness in his decision to die by providing lethal injections.

While ethics and law are concerned with different concepts of right and wrong, in medicine they find common ground in their fundamental principles. Both law and ethics in medicine rest on the principle of self-determination by competent individuals, beneficence (or at least nonmaleficence) on the part of medical practitioners, and a concept of justice as fairness to be afforded to all patients by both medical practitioners and society.

Law and medical practice

The law indirectly influences medical practice by structuring the delivery and financing of medical services and it does so directly in three major ways: licensure requirements; restrictions on practice; and redress for wronged patients.

GOVERNMENT FINANCING

Societies differ as to how much medical care is considered adequate. In the United States, the President's Commission for the Study of Ethical Problems in Medicine concluded in 1983 that even though access to health care might not be properly considered a legal right, "society has a moral obligation to ensure that everyone has access to adequate care without being subject to excessive burdens." Programs for governmental payment of medical bills are an example of law used as an instrument to effect social policy. Changing the laws regarding financing rules may, in fact, have more of an impact on medical practice than any other legal interventions.

LICENSE REQUIREMENTS FOR PRACTICE

Legislative approaches to health-care licensing have been divided into three general categories by the legal scholar Jan Stepan: (1) exclusive or monopolistic systems, in which only the practice of modern, scientific medicine by professionals is lawful; (2) tolerant systems, in which only scientific medicine is officially sanctioned, although traditional medicine is tolerated; and (3) inclusive or integrated systems, in which traditional medicine is either recognized as a special part of the health system, or the integration of two or more systems of health care is officially promoted.

The exclusive systems. The exclusive system of licensing, which first emerged in Europe, involves granting licenses only to those who meet certain formal educational requirements; they are often required to pass a standardized examination and to demonstrate that they are of good moral character. The educational requirements usually apply to physicians, nurses, pharmacists, dentists, and other allied health professionals. Most medical schools in

the world have similar courses of study, and a person who is graduated from almost any of them receives the M.D. degree (doctor of medicine) and is entitled to sit for a licensure examination. In Britain and many Commonwealth countries the graduate receives a bachelor of medicine and a bachelor of surgery or bachelor of chirurgery (M.B., B.S.; or M.B., Ch.B.); the M.D. is a higher degree that requires further study. The license to practice, called registration, is granted by the General Medical Council, which was established by the Medical Act of 1858.

Because of the similarities between educational and practice patterns, physicians often may be allowed to practice in countries other than those in which they are registered. In 1975 the Council of the European Economic Communities, for example, issued directives for the national recognition of primary and specialist medical qualifications so as to facilitate the free movement of doctors among the member states of the European Economic Community (EEC). Steps have been taken in most EEC countries to implement these directives. An Advisory Committee in Medical Training has been established by the EEC to assist in the formation of a common standard.

Legislatures originally adopted exclusive licensing in an effort to protect the public from untrained persons on the belief that the public was not qualified to judge medical skills and could easily be taken advantage of by unqualified practitioners. Licenses are generally granted for life and can be lost only by criminal conduct, gross incompetence, or mental disability. Virtually every country in the world, with the exceptions of China and Japan, initially followed the European model of exclusive licensing, and the type of health care this model imposed is based on Western science. What may have been reasonable for Europe and North America, however, makes little sense in rural Africa or South America, and it has become clear that more flexibility is needed to ensure better access to health care in many of the developing countries.

Exclusive approaches to medical care exist in almost every state in the United States and in France, Belgium, and most other European countries. In the United States, for example, it is a criminal offense for anyone to practice medicine without a license. The practice of medicine is usually defined to include diagnosis, treatment, prescription, and surgery. Physicians may lawfully practice in any branch of medicine and any of its subspecialties. Further education, training, and certification as a specialist is available through private, voluntary organizations. All other licensed health professionals, however, are strictly limited by statute to the areas in which they may practice. For example, dentists must only do work relating to dentistry.

Because the U.S. states have retained direct powers in health care, all licenses for health professionals are granted by the individual states. Licensing is regulated by administrative agencies, usually called state licensing boards, and licenses are valid only in the state issuing the license. Most states, however, grant a license by reciprocity upon request from a licensed physician, nurse, pharmacist, or dentist who is in good standing in another state. Licensing rules governing dentistry and pharmacy are parallel to medical licensure: the candidate must have been graduated from a professional school, have passed an examination, and be licensed by a state board set up for that purpose.

Registered nurses (R.N.'s) are required to have been graduated from a preparatory nursing program (with a bachelor's degree, associate degree, or diploma from a hospital-based program) and to have passed an examination. Licensed practical nurses (L.P.N.'s rather than R.N.'s), on the other hand, generally have to complete a one-year vocational program after graduating from high school. In almost all countries, national and local nursing associations are instrumental in proposing legislation for licensing or registration and in securing amendments to broaden the scope of nursing practices.

The tolerant systems. Tolerant systems operate in Britain, Germany, the Scandinavian countries, and some U.S. states, such as California. In Britain, for example, there is no formal legal monopoly of medicine. The protected status is that of registered medical practitioner. Under the Medical Act of 1978, persons who have fulfilled

Licensing

statutory education and examination requirements are entitled to be registered. Registered physicians have certain exclusive rights, such as employment by the National Health Service, prescribing, issuing medical certificates, and holding appointments in public hospitals. It is an offense to imply falsely that one is registered, but it is not an offense to otherwise engage in healing. Osteopaths, chiropractors, and acupuncturists may practice in the private sector. The Dental Act of 1984, however, gives dentists a monopoly by making it an offense to practice dentistry when not registered as a dentist. Sweden is similar in that anyone may be a healer, but it is a crime to claim to be a physician when not so recognized by the state.

The inclusive and integrated systems. The inclusive and integrated systems are represented by countries such as India and China. In India the individual states regulate the practice of the health professions. A Medical Council maintains the Indian Medical Register, and any person on this register may practice anywhere in India. There are also state medical registers for the individual states, and people in these registers are also in the Indian Medical Register. Those who have obtained a degree from a university, usually M.B., B.S., or the qualification of licensed medical practitioner (L.M.P.) are entered in these registers. Qualifications from many foreign countries are also recognized for entry to the registers. Although it is now a crime to practice medicine without a license in India, this law was only enacted in 1970, and it does not apply to anyone who had practiced medicine without a license for at least five years before 1970.

China provides an example of an integrated system in which traditional Chinese doctors and Western-trained physicians practice side by side. The country has also developed a system in which certain persons receive a short, intensive training course and then provide primary health care in areas otherwise not served by health-care workers.

LEGAL RESTRICTIONS ON PRACTICE

Determination of death. The law generally supports customary medical practice and provides the medical profession with a great deal of autonomy. A dramatic example is the determination of death and the issuance of a death certificate. In almost every country of the world a physician declares a person dead and issues a death certificate after a determination of death is made in accordance with accepted medical standards. A question that recently appeared was whether physicians should continue to be given the authority to declare a person dead if the medical profession were to adopt whole brain death as an acceptable definition of death (instead of the past definition of irreversible cessation of respiration and heartbeat). A mechanical respirator can artificially maintain the respiration and circulation of a person whose functions would cease without such mechanical support. In the late 1960s the potentials of organ transplantation from such persons were becoming realized, and the seeming futility of devoting limited medical resources to maintaining circulation under such circumstances was of growing concern. Physicians began proposing that irreversible cessation of brain activity be used as an alternative definition of death.

Since that time, most Western countries have adopted this definition, by either continuing to permit physicians to declare death, passing a specific statute endorsing this definition, or issuing court opinions giving approval to physicians' declarations of death in such circumstances. The law, in short, has continued to defer to medical practice in the definition of death itself.

The countries that have not adopted brain death criteria have not done so primarily for cultural and religious reasons. For example, Japan has refused to adopt a brain-based definition of death in part because it would conflict with religious tenets that require the death of all major organs prior to a pronouncement of death. Accordingly, such medical techniques as heart transplantation cannot be performed in Japan (or in any country that does not accept a brain-death definition). Harvesting a beating heart from a person is considered in such countries to be homicide, even though brain activity has ceased and respiration is being maintained artificially.

Termination of pregnancy. Physicians have broad legal authority and discretion in other controversial areas as well. In countries like the United States, for example, where the law permits termination of pregnancy prior to fetal viability, it is for the physician, and not the state, to determine whether or not an individual fetus is viable (*i.e.*, capable of living independently of its mother); the determination must be made consistent with accepted medical criteria. When the Supreme Court declared that the decision to terminate a pregnancy was protected by the U.S. Constitution, it emphasized that the decision should be made by a woman and her physician. The law may restrict procedures like termination of pregnancy, sterilization, and even birth control to mature minors and adults; but it is generally left to the physician to determine if a patient is mature or competent to consent.

Sudden death. Unexpected deaths must be reported to public authorities so that a determination can be made as to whether the death was a homicide, a suicide, or an accident. It is up to an investigator, whether he is called a medical examiner or coroner, or has some other title, to make a preliminary finding and then to refer the case to the police or prosecuting authorities if criminal activity is suspected. In most Western countries this person has either legal training, medical training, or both. In the United States, some jurisdictions use an elected coroner (who may or may not have medical training); however, the trend is toward a medical examiner system in which a physician is appointed to the post. Coroners in London are qualified both in medicine and in law. The principal evidence received by the coroner is the report of the autopsy carried out on the body of the deceased.

Public reporting. Physicians may also be expected to report certain patients or occurrences to public authorities. For example, some communicable diseases are required to be reported to public health officials. Suspected child abuse and gunshot wounds may have to be reported to an authority (such as the child welfare authority or the police). Public reporting tends to put the physician in the position of being an agent of the state rather than of his patient, making mandatory reporting an uncomfortable duty. It is not surprising, therefore, that the incidence of public reporting is often thought to be lower than may be warranted.

In extreme cases, physicians may also have a duty to warn specific individuals who may be at serious or mortal risk from their patients. For example, the California Supreme Court decreed that a psychologist had a "duty to warn" a person whom his patient had threatened to kill, if the psychologist believed "or should have believed" the threat to be real. In that particular case the patient, a graduate student, left the psychologist's care and murdered his former girlfriend, an action that the psychologist believed was so likely that he had sought unsuccessfully to have his patient involuntarily committed.

LEGAL REDRESS

When patients are injured by medical negligence the remedies they can pursue depend upon the country's legal system. In the United States, for example, lawsuits against physicians for negligent injury are not considered unusual.

Malpractice, or professional negligence, is the failure of a health-care provider (for example, a physician, dentist, nurse, or pharmacist) to exercise the ordinary care and skill a reasonably prudent, qualified person would exercise under the same or similar circumstances. The practitioner does not guarantee the outcome but must use diligence and ordinary skill in the treatment of a patient.

A valid malpractice claim must have four elements: duty, breach, damages, and causation. The practitioner must be shown to have a relationship to the patient (which establishes a duty to exercise ordinary care); must have breached that duty (as measured by the applicable standard of care); and through the breach must have caused the patient physical and monetary damages.

The central concern for physicians is usually to establish the standard of care through expert testimony, which may simply be the testimony of another qualified physician. Such testimony is necessary because the standard of med-

Effect of public reporting

ical practice is not something a lay jury is familiar with. Expert witnesses may themselves rely on the standards that have been set down by one or more of the medical speciality organizations, such as the American College of Obstetricians and Gynecologists. These medical speciality organizations provide certification to physicians who have fulfilled postgraduate training and practice requirements in the speciality. They maintain the standards necessary to practice in the specialities, and they provide reasonable assurance to patients that these standards will be upheld. Nonconformance with such standards by a specialist is evidence of negligence, although it is not conclusively negligence (the practitioner may have a valid excuse for not following custom, such as an emergency situation or lack of equipment). Conformance with the standards is evidence of due care, but it is not conclusive because other factors may have caused the physician's action to be imprudent under the circumstances.

If a practitioner consistently performs below the profession's standard of care (*i.e.*, the practitioner is a negligent physician who does not actually harm anyone) the remedy is not a malpractice action but a complaint to the licensing or registration authority to have the individual disciplined. Disciplinary action by public licensing authorities, however, is unusual.

Functions
of
malpractice

Medical malpractice actions have three basic functions: quality control, compensation for harm, and emotional vindication, all of which are achieved to varying degrees. Quality control is probably best achieved, since the standard of care is set by physicians themselves and enforced by patients and juries. Compensation for harm, on the other hand, is skewed toward major injuries. Attorneys in the United States, for example, represent malpractice cases on a contingency fee basis—*i.e.*, they are paid a proportion (usually 20–40 percent) of the total amount awarded to the plaintiff. Patients who suffer less severe injuries may have little redress for compensation. In countries that have a system of national health insurance, compensation for harm may not be a major issue (since all medical bills are paid regardless of cause). And countries with comprehensive social services for all citizens, like Sweden and New Zealand, have effectively developed “no fault” compensation systems. But in the United States, where more than 35,000,000 people do not have any form of health-care insurance, lack of coverage can transform a medically induced injury into a financial catastrophe.

Emotional vindication is a measure of the consumer's ability to make a complaint as well as to get a satisfactory response. A comparison in consumer complaints between the United States and Britain shows significant differences. In 1981 an estimated 700 writs for medical malpractice were issued in Britain. In 1983 about 42,000 claims were filed in the United States. These figures yield a per capita difference (even when the difference in years and between claims and writs is considered) that shows that Americans file claims against physicians more than 10 times as often as their British counterparts. The law professor Frances Miller has noted that many cultural and practical reasons serve to explain this difference, including different legal systems and rules, access to attorneys and courts, the method of paying for medical expenses, the special status of the National Health Service in Britain, and the existence of alternative complaint procedures.

Patients' rights

In addition to granting patients means for the effective redress for negligent injury (which increases the cost of malpractice insurance for physicians—and thus the cost of medical care), malpractice litigation has also promoted what have come to be called patients' rights.

Funda-
mental
premises

Patients' rights are based upon two fundamental premises: (1) the patient has certain interests, many of which may be properly described as rights, that are not automatically forfeited by entering into a relationship with a doctor or health-care facility; and (2) doctors and health-care facilities may fail to recognize the existence of these interests and rights, fail to provide for their protection or assertion, and frequently limit their exercise without recourse.

Perhaps the most important development in patient rights has been that in the United States regarding the doctrine of informed consent. This doctrine requires physicians to share certain information with patients before asking for their consent to treatment. The doctrine is particularly applicable to the use of surgery, drugs, and invasive diagnostic procedures that carry risks. It requires the physician to describe the procedure or treatment recommended and to list its major risks, benefits, alternatives, and likely prospect for recuperation. The purpose is to promote self-determination by patients on the theory that it is the patient who has the most at stake in treatment and who relies largely on the physician for such information. British courts have rejected this formulation on the basis that the average British citizen does not want such information, and British physicians do not generally provide it unless requested.

In 1972 the American Hospital Association adopted a patient bill of rights based on the premise that “[the] traditional physician-patient relationship takes on a new dimension when care is rendered within an organizational structure . . . the institution itself also has a responsibility to the patient.” The text of the American Hospital Association bill of patient rights calls for the rights of the patient to respectful care, complete medical information, information necessary for informed consent, refusal of treatment, privacy, confidentiality, response to requests for service, information on other institutions involved in the patient's care, refusal of participation in research projects, continuity of care, examination and explanation of financial charges, and knowledge of hospital regulations.

In the mid-1970s the Parliamentary Assembly of the Council of Europe submitted a draft recommendation to its member nations suggesting that all necessary action be taken to ensure that the sick can receive relief from their suffering and that patients be adequately prepared for death. They further recommended that commissions be established to study the issue of euthanasia and that physicians acknowledge that the sick have a right to full information, if requested, on their illness and the proposed treatment. The council listed the following basic rights for the humane and dignified treatment of patients: the right to freedom, to personal dignity and integrity, to information, to proper care, and to not suffer.

No one bill of rights is suitable for every health care institution, and specialty hospitals, such as maternity and pediatric hospitals, may require different approaches. Among the basic rights of a patient should be the right to clear communication: accurate information concerning possible medical care and procedures; informed participation in all decisions involving the patient's health-care program; and a clear, concise explanation of all proposed procedures, including possible risks, side effects, and problems related to recuperation.

Patients also should have rights regarding quality of care: a right to an accurate evaluation of their condition and prognosis without treatment; knowledge of the identity and professional status of those providing services; information contained in their medical record; access to consultant specialists; and refusal of treatment.

The patient should have basic human rights: the right to privacy of both person and information; of access to people outside the health-care facility; and to leave the health-care facility regardless of his condition.

The goals of such a system are to protect patients, especially those at a disadvantage within the health-care system (*e.g.*, the young, the illiterate, the uncommunicative, and those without relatives); to make available to those who seek it the opportunity to participate actively with the physician as a partner in a personal health-care program; and to put into proper perspective medical technology and pharmaceutical advances.

Until the 1960s law and medicine met only in the courtroom, and then usually only in cases involving pathology or psychiatry. Since then, however, civil litigation, public financing, and ethical issues have grown, at least partially as a result of the incredible successes of medicine. These successes have increased public expectations and increased the cost of medicine; they also have made decisions about

Proposed
basic rights

terminating care more ambiguous. Enhancing patients' rights is one modern concern on which both medical and legal practitioners agree. (G.J.A.)

BIBLIOGRAPHY

General works: Webster's Medical Desk Dictionary (1986), is a reference source for the layman. *The Oxford Companion to Medicine*, 2 vol., edited by JOHN WALTON, PAUL B. BEESON, and RONALD BODLEY SCOTT (1986), is a comprehensive text of 20th-century developments and persons.

History: The literature on the history of medicine covers all topics and periods and includes biographies as well as descriptions of the development of hospitals, research institutes, health care, and medical education in different countries. Introductory studies include GEORGE T. BETTANY, *Eminent Doctors: Their Lives and Their Work*, 2 vol. (1885, reprinted 1972); ARTURO CASTIGLIONI, *A History of Medicine*, 2nd rev. ed. (1947; originally published in Italian, 1927), a classic work; FIELDING H. GARRISON, *An Introduction to the History of Medicine*, 4th rev. ed. (1929, reprinted 1967), a scholarly history; DOUGLAS GUTHRIE, *A History of Medicine*, rev. ed. (1958); HOWARD W. HAGGARD, *Devils, Drugs, and Doctors: The Story of the Science of Healing from Medicine-Man to Doctor* (1929, reprinted 1980); RICHARD H. MEADE, *An Introduction to the History of General Surgery* (1968), a well-documented work on developments in surgery on separate organs; CHARLES SINGER and E. ASHWORTH UNDERWOOD, *A Short History of Medicine*, 2nd ed. (1962); and PHILIP RHODES, *An Outline History of Medicine* (1985).

Ancient traditions of non-Western medicine are presented in P. KUTUMBIAH, *Ancient Indian Medicine* (1962); HEINRICH R. ZIMMER, *Hindu Medicine* (1948, reprinted 1979); EDWARD H. HUME, *The Chinese Way in Medicine* (1940, reprinted 1975); PAUL U. UNSCHULD, *Medicine in China: A History of Ideas* (1985; originally published in German, 1980); and EDWARD G. BROWNE, *Arabian Medicine* (1921, reprinted 1983).

For developments from the origin of Western medicine to the end of the 18th century, see WILLIAM G. BLACK, *Folk-Medicine: A Chapter in the History of Culture* (1883, reprinted 1970); W.H.R. RIVERS, *Medicine, Magic, and Religion* (1924, reprinted 1979), a comprehensive treatment of primitive medicine; JOHN SCARBOROUGH, *Roman Medicine* (1969, reprinted 1976); ROBERT S. GOTTFRIED, *Doctors and Medicine in Medieval England, 1340-1530* (1986); A. WEAR, R.K. FRENCH, and I.M. LONIE (eds.), *The Medical Renaissance of the Sixteenth Century* (1985); KATHARINE PARK, *Doctors and Medicine in Early Renaissance Florence* (1985); and GUY WILLIAMS, *The Age of Agony: The Art of Healing, c. 1700-1800* (1975, reprinted 1986).

Medicine and surgery during the 19th and 20th centuries are the subject of CARL J. PFEIFFER, *The Art and Practice of Western Medicine in the Early Nineteenth Century* (1985); THOMAS E. KEYS, *The History of Surgical Anesthesia*, rev. ed. (1963, reprinted 1978); M.H. ARMSTRONG DAVISON, *The Evolution of Anesthesia* (1965); ROBERT G. RICHARDSON, *The Scalpel and the Heart* (1970; U.K. title, *The Surgeon's Heart: A History of Cardiac Surgery*, 1969); JOHN S. HALLER, JR., *American Medicine in Transition, 1840-1910* (1981); RUTH J. ABRAM (ed.), *Send Us a Lady Physician: Women Doctors in America, 1835-1920* (1985); A. MCGEEHE HARVEY, *Science at the Bedside: Clinical Research in American Medicine, 1905-1945* (1981), a discussion of the institutionalization of clinical research; and LAWRENCE GALTON, *Med Tech: The Layperson's Guide to Today's Medical Miracles* (1985), a historical dictionary.

Medical practice: GEORGE ROSEN, *The Structure of American Medical Practice, 1875-1941* (1983), is a historical study. Particular kinds of medical practice are explored in WESLEY FABB and JOHN FRY (eds.), *Principles of Practice Management in Primary Care* (1984); SIR DOUGLAS BLACK et al., *Inequalities in Health: The Black Report*, edited by PETER TOWNSEND and NICK DAVIDSON (1982); DAVID SANDERS and RICHARD CARVER, *The Struggle for Health: Medicine and the Politics of Underdevelopment* (1985); and V. DJUKANOVIC and E.P. MACH (eds.), *Alternative Approaches to Meeting Basic Health Needs in Developing Countries: A Joint UNICEF/WHO Study* (1975). Also see the articles of such journals as *Private Practice* (monthly) and *Modern Healthcare* (semimonthly). For a view of alternative medicine, see DOUGLAS STALKER and CLARK GLYMOUR (eds.), *Examining Holistic Medicine* (1985); and RICHARD GROSSMAN, *The Other Medicines* (1985). The variety of roles in the health-care profession are the subject of LOUISE SIMMERS, *Diversified Health Occupations* (1983); C. WESLEY EISELE, WILLIAM R. FIFER, and TOMA C. WILSON, *The Medical Staff and the Modern Hospital* (1985); and ELI GINZBERG (ed.), *From Physician Shortage to Patient Shortage: The Uncertain Future of Medical Practice* (1986).

Among the many books devoted to the subject of medical and health-care education are the following historical discussions: ABRAHAM FLEXNER, *Medical Education in the United States and Canada: A Report to the Carnegie Foundation for the Ad-*

vancement of Teaching (1910, reprinted 1973); and KENNETH M. LUDMERER, *Learning to Heal: The Development of American Medical Education* (1985). For special information, see the following official publications: ASSOCIATION OF AMERICAN MEDICAL COLLEGES, *AAMC Directory of American Medical Education, 1986-87*, 33rd ed. (1986); *Medical School Admission Requirements, 1988-89*, 38th ed. (1987); and *Physicians for the Twenty-first Century: Report to the Project Panel on the General Professional Education of the Physician and College Preparation for Medicine* (1984). Studies include MOHAN L. GARG and WARREN M. KLEINBERG, *Clinical Training and Health Care Costs: A Basic Curriculum for Medical Education* (1985); and MARJORIE PRICE WILSON and CURTIS P. MCLAUGHLIN, *Leadership and Management in Academic Medicine* (1984).

For new developments in medical education, see the periodicals *The Journal of Medical Education* (monthly), *Medical Education* (bimonthly), and *WHO Chronicle* (bimonthly). Opportunities for continuing medical education appear semiannually in *JAMA: The Journal of the American Medical Association* (weekly). Education for health-care professionals other than physicians is the subject of COMMITTEE FOR THE STUDY OF NURSING EDUCATION (U.S.), *Nursing and Nursing Education in the United States* (1923, reprinted 1984), a classic yet still timely report on nursing training; BRYN DAVIS (ed.), *Nursing Education: Research and Developments* (1987); SHIRLEY A. BADASCH and DOREEN S. CHESEBRO, *The Health Care Worker: An Introduction to Health Occupations* (1985); JOSEPH S. GREEN et al. (eds.), *Continuing Education for the Health Professions* (1984); and T. DONALD RUCKER (ed.), *Pharmacy: Career Planning and Professional Opportunity* (1981).

Medical institutions: For historical discussions of hospitals, see TIMOTHY S. MILLER, *The Birth of the Hospital in the Byzantine Empire* (1985); GUENTER B. RISSE, *Hospital Life in Enlightenment Scotland: Care and Teaching at the Royal Infirmary of Edinburgh* (1986); and RICHARD H. THURM, *For the Relief of the Sick and Disabled: The U.S. Public Health Service Hospital in Boston, 1799-1969* (1972). Organization and administration of modern hospitals are discussed in I. DONALD SNOOK, JR., *Hospitals: What They Are and How They Work* (1981); JONATHAN S. RAKICH and KURT DARR (eds.), *Hospital Organization and Management*, 3rd ed. (1983); THOMAS CHOI, ROBERT F. ALLISON, and FRED MUNSON, *Governing University Hospitals in a Changing Environment* (1986); HOWARD J. BERMAN, LEWIS E. WEEKS, and STEVEN F. KUKLA, *The Financial Management of Hospitals*, 6th ed. (1986); BRADFORD H. GRAY (ed.), *The New Health Care for Profit: Doctors and Hospitals in a Competitive Environment* (1983); and EVERETT A. JOHNSON and RICHARD L. JOHNSON, *Hospitals Under Fire: Strategies for Survival* (1986).

For developments in the field of hospitals, see the periodical *Hospitals* (semimonthly), published by the American Hospital Association. Social and psychological aspects of hospital life are the subject of GEOFFREY C. ROBINSON and HEATHER F. CLARKE, *The Hospital Care of Children: A Review of Contemporary Issues* (1980); and JUDITH WILSON ROSS, *Handbook for Hospital Ethics Committees* (1986). Hospital building is analyzed in W. PAUL JAMES and WILLIAM TATTON-BROWN, *Hospitals: Design and Development* (1986); OWEN B. HARDY and LAWRENCE P. LAMMERS, *Hospitals, the Planning and Design Process*, 2nd ed. (1986); and I. DONALD SNOOK, JR., and KATHIRYN M. RUCK (eds.), *Using Hospital Space Profitably* (1987).

For organization and administration of other health facilities, see MARTIN P. CHARNS and MARGUERITE J. SCHAEFFER, *Health Care Organizations, a Model for Management* (1983); and BRADFORD H. GRAY (ed.), *For-Profit Enterprise in Health Care* (1986). The following works discuss aspects of the activities of outpatient facilities: DREXEL TOLAND and SUSAN STRONG, *Hospital-Based Medical Office Buildings*, 2nd ed. (1986); PETER M. FRIEND and JOHN M. SILVER (eds.), *Freestanding Emergency Centers: A Guide to Planning, Organization, and Management* (1985); and ROBERT KOHN, *The Health Centre Concept in Primary Health Care* (1983).

Facilities for geriatric care are discussed in WILLIAM B. BRIGER and WILLIAM R. POMERANZ, *Nursing Home Development: A Guide to Planning, Financing, and Constructing Long-Term Care Facilities* (1985); and COLLEEN L. JOHNSON and LEISLIE A. GRANT, *The Nursing Home in American Society* (1985). Medical institutions for terminal care are the subject of CHARLES A. CORR and DONNA M. CORR (eds.), *Hospice Care: Principles and Practice* (1983); GLEN W. DAVIDSON (ed.), *The Hospice: Development and Administration*, 2nd ed. (1985); JACK MCKAY ZIMMERMAN, *Complete Care for the Terminally Ill*, 2nd ed. (1986); and CHARLES A. CORR and DONNA M. CORR (eds.), *Hospice Approaches to Pediatric Care* (1985); as well as of articles in *The American Journal of Hospice Care* (bimonthly).

Public health: FRASER BROCKINGTON, *World Health*, 3rd ed. (1975), is a comprehensive discussion of public health concepts and the World Health Organization. JOHN BRYANT, *Health &*

the Developing World (1969, reprinted 1972), studies in health-care in Africa, Asia, and Latin America. JOHN M. LAST (ed.), *Public Health and Preventive Medicine*, 11th ed. (1980), is a definitive text. Later surveys of the organized effort to protect and improve community health include DEREK ERASER, *The Evolution of the British Welfare State: A History of Social Policy Since the Industrial Revolution*, 2nd ed. (1984); ROBERT LANZA (ed.), *Medical Science and the Advancement of World Health* (1985); and GRACE BUDRYS, *Planning for the Nation's Health: A Study of Twentieth-Century Developments in the United States* (1987).

Nursing: The long historical tradition of the nursing occupation is studied in JOHN HOWARD, *An Account of the Principal Lazarettos in Europe* (1789); BRIAN ABEL-SMITH, *A History of the Nursing Profession* (1960, reprinted 1975); VERN L. BULLOUGH and BONNIE BULLOUGH, *The Emergence of Modern Nursing*, 2nd ed. (1969); CECIL BLANCHE WOODHAM SMITH, *Florence Nightingale, 1820-1910* (1950, reissued 1983); MONICA E. BALLY, *Florence Nightingale and the Nursing Legacy* (1986); BARBARA MELOSH, *The Physician's Hand: Work Culture and Conflict in American Nursing* (1982); and M. PATRICIA DONAHUE, *Nursing, the Finest Art: An Illustrated History* (1985).

ESTHER LUCILE BROWN, *Nursing Reconsidered: A Study of Change*, 2 vol. (1970-71), reviews the developments of the 1960s and analyzes new functions performed by nurses. See also MONICA E. BALLY, *Nursing and Social Change*, 2nd ed. (1980); BONNIE BULLOUGH, VERN L. BULLOUGH, and MARY CLAIRE SOUKUP, *Nursing Issues and Nursing Strategies for the Eighties* (1983); and NORMA L. CHASKA (ed.), *The Nursing Profession: A Time to Speak* (1983). Nursing practices are explored in SANDRA DEBELLA, LEONIDE MARTIN, and SANDRA SIDDALL, *Nurses' Role in Health Care Planning* (1986); MARY H. BROWNING (comp.), *The Nursing Process in Practice* (1974), and *Nursing and the Aging Patient* (1974); MARY H. BROWNING and EDITH P. LEWIS (comps.), *Nursing and the Cancer Patient* (1973), and *The Nurse in Community Mental Health* (1972); ANDREW JAMETON, *Nursing Practice: The Ethics Issues* (1984); and CAROL REN KNEISL and SUEANN WOOSTER AMES (eds.), *Adult Health Nursing: A Biopsychosocial Approach* (1986).

Dentistry: PIERRE EAUCHARD, *The Surgeon Dentist: or, Treatise on the Teeth*, 2 vol. in 1 (1946, reprinted 1969; originally published in French, 2nd ed., 1746) is a classic work. See also VINCENZO GUERINI, *A History of Dentistry: From the Most Ancient Times Until the End of the Eighteenth Century* (1909, reprinted 1967); and BERNHARD WOLE WEINBERGER, *An Introduction to the History of Dentistry, with Morality & Dental Chronology & Bibliographic Data*, 2 vol. (1948, reprinted 1981). MALVIN E. RING, *Dentistry: An Illustrated History* (1985), is a readable comprehensive survey. For a broad range of current information on dentistry, see *Dental Abstracts* (monthly), a compilation of abstracts of articles from periodicals; and *The Year Book of Dentistry*, an annual review of selected articles.

R.A. CAWSON and J.W. EVESON, *Oral Pathology and Diagnosis* (1987), is a text and colour atlas. Oral and dental surgery is the subject of W. HARRY ARCHER, *Oral and Maxillofacial Surgery*, 5th ed., 2 vol. (1975); GUSTAV O. KRUGER, *Textbook of Oral and Maxillofacial Surgery*, 6th ed. (1984); and DANIEL M. LASKIN (ed.), *Oral Surgery* (1985). Other dental specialties are covered in JOHN M. DAVIS, DAVID B. LAW, and THOMPSON M. LEWIS, *An Atlas of Pedodontics*, 2nd ed. (1981); RALPH E. McDONALD and DAVID R. AVERY (eds.), *Dentistry for the Child and Adolescent*, 4th ed. (1983); THOMAS M. GRABER and BRAINERD E. SWAIN, *Orthodontics: Current Principles and Techniques* (1985); LEE W. GRABER (ed.), *Orthodontics: State of the Art, Essence of the Science* (1986); LOUIS I. GROSSMAN, *Endodontic Practice*, 10th ed. (1981); IRVING GLICKMAN, *Glickman's Clinical Periodontology*, 6th ed., edited by EERMIN A. CARRANZA (1984); HENRY M. GOLDMAN and D. WALTER COHEN, *Periodontal Therapy*, 6th ed. (1980); HENRY M. GOLDMAN et al. (eds.), *Current Therapy in Dentistry*, 7th ed. (1980); JOHN E. JOHNSTON, *Johnston's Modern Practice in Fixed Prosthodontics*, 4th ed., edited by ROLAND W. DYKEMA, CHARLES J. GOODACRE, and RALPH W. PHILLIPS (1986); EUGENE WILLIAM SKINNER, *Skinner's Science of Dental Materials*, 8th ed., edited by RALPH W. PHILLIPS (1982); ROBERT G. CRAIG (ed.), *Restorative Dental Materials*, 7th ed. (1985); and ERANK M. MCCARTHY (ed.), *Emergencies in Dental Practice: Prevention and Treatment*, 3rd ed. (1979).

Prophylactic measures in dentistry are discussed in JOSEPH L. BERNIER and JOSEPH C. MUHLER (eds.), *Improving Dental Practice Through Preventive Measures*, 3rd ed. (1975); and ABRAHAM E. NIZEL, *Nutrition in Preventive Dentistry: Science and Practice*, 2nd ed. (1981). DAVID E. STRIEFLER, WESLEY O. YOUNG, and BRIAN A. BURT, *Dentistry, Dental Practice, & the*

Community, 3rd rev. ed. (1983); and JAMES MORSE DUNNING, *Principles of Dental Public Health*, 4th ed. (1986), are studies of the relationship between the dentist and society.

Pharmacy: EDWARD KREMERS, *Kremers and Urdang's History of Pharmacy*, 4th ed., rev. by GLENN SONNEDECKER (1976), provides comprehensive coverage of the topic. Useful information can be found in UNITED STATES PHARMACOPEIAL CONVENTION, *Advice for the Patient: Drug Information in Lay Language*, 6th ed. (1986). Modern studies of practices in pharmacy include WILLIAM E. HASSAN, JR., *Hospital Pharmacy*, 5th ed. (1986); THOMAS R. BROWN and MICKEY C. SMITH (eds.), *Handbook of Institutional Pharmacy Practice*, 2nd ed. (1986); PAUL TURNER, ALAN RICHENS, and PHILIP ROUTLEDGE, *Clinical Pharmacology*, 5th ed. (1986); JOSEPH P. REMINGTON, *Remington's Pharmaceutical Sciences*, 17th ed., edited by ALEONSO R. GENNARO et al. (1985); D.R. LAURENCE and P.N. BENNETT, *Clinical Pharmacology*, 5th ed. (1980); and ALERED GOODMAN GILMAN et al. (eds.), *Goodman and Gilman's The Pharmacological Basis of Therapeutics*, 7th ed. (1985).

Legal medicine: H. TRISTRAM ENGELHARDT, JR., *The Foundations of Bioethics* (1986), gives an overview of medical ethics. EMILY ERIEDMAN (ed.), *Making Choices: Ethics Issues for Health Care Professionals* (1986), is a selection of articles on decision-making in the rationing of health care and in the provision of critical care. UNITED STATES. PRESIDENT'S COMMISSION FOR THE STUDY OF ETHICAL PROBLEMS IN MEDICINE AND BIOMEDICAL AND BEHAVIORAL RESEARCH, *Securing Access to Health Care: A Report on the Ethical Implications of Differences in the Availability of Health Services*, 3 vol. (1983), is an authoritative document.

Issues accompanying the medical problems of special neonatal care are discussed in EARL E. SHELF, *Born to Die?: Deciding the Fate of Critically Ill Newborns* (1986); and SHERMAN ELIAS and GEORGE J. ANNAS, *Reproductive Genetics and the Law* (1987). The ethics of general critical care are the subject of STUART J. YOUNGNER (ed.), *Human Values in Critical Care Medicine* (1986); and ALONZO L. PLOUGH, *Borrowed Time: Artificial Organs and the Politics of Extending Lives* (1986). Patients' rights in view of the law of informed consent are explored in GEORGE J. ANNAS, *The Rights of Hospital Patients: The Basic ACLU Guide to a Hospital Patient's Rights* (1975); JOANNE LYNN (ed.), *By No Extraordinary Means: The Choice to Forgo Life-Sustaining Food and Water* (1986); WILLIAM J. WINSLADE and JUDITH WILSON ROSS, *Choosing Life or Death: A Guide for Patients, Families, and Professionals* (1986); and JAMES RACHELS, *The End of Life: Euthanasia and Morality* (1986).

Laws and legislation regulating health care are analyzed in GEORGE J. ANNAS, LEONARD H. GLANTZ, and BARBARA E. KATZ, *The Rights of Doctors, Nurses, and Allied Health Professionals: A Health Law Primer* (1981); SAMUEL C. URSU (ed.), *Symposium on Legal Considerations in Dentistry* (1982); ERANCES H. MILLER, "Medical Malpractice Litigation: Do the British Have a Better Remedy?," *American Journal of Law & Medicine*, 11:433-463 (1986), a comparison of the systems of medical malpractice litigation in the United States and Britain; PATRICIA M. DANZON, *Medical Malpractice: Theory, Evidence, and Public Policy* (1985); RAYMOND G. DEVRIES, *Regulating Birth: Midwives, Medicine, & the Law* (1985); and ROBERT D. MILLER, *Problems in Hospital Law*, 5th ed. (1986).

A summary of the major models of licensing of health-care professionals throughout the world is given in JAN STEPAN, "Traditional and Alternative Systems of Medicine: A Comparative Review of Legislation," *International Digest of Health Legislation*, 36:281-341 (1985). On professional accreditation and licensing, see DONALD G. LANGSLEY (ed.), *Legal Aspects of Certification and Accreditation* (1983); CLARE LABAR, *Statutory Requirements for Licensure of Nurses* (1985); and COUNCIL OF STATE GOVERNMENTS, *State Credentialing of the Health Occupations and Professions* (1986).

Forensic medicine is the subject of R.B.H. GRADWOHL, *Gradwohl's Legal Medicine*, 3rd ed., edited by FRANCIS E. CAMPS (1976), a comprehensive, fully illustrated text of traditional legal medicine. JOHN GLAISTER, *Glaister's Medical Jurisprudence and Toxicology*, 13th ed., edited by EDGAR RENTOU and HAMILTON SMITH (1973), is a comprehensive text on forensic pathology and toxicology. See also WERNER U. SPITZ and RUSSELL S. FISHER (eds.), *Medicolegal Investigation of Death: Guidelines for the Application of Pathology to Crime Investigation*, 2nd ed. (1980); SIR JOHN JERVIS, *Jervis on the Office and Duties of Coroners*, 10th ed., edited by PAUL MATTHEWS and J.C. EOREMAN (1986); and WILLIAM J. CURRAN and E. DONALD SHAPIRO, *Law, Medicine, and Forensic Science*, 3rd ed. (1982).

(P.Rh./Fr.H./G.J.A./Ha.Sc./Ed.)

Melbourne

Although the central city is the home of fewer than 100,000 people, Melbourne is the core of an extensive metropolitan area and is the world's most southerly urban area with a population of more than 1,000,000. It is situated on the southeastern coast of Australia, at the head of Port Phillip Bay, and is the capital of the state of Victoria. In Australia its metropolitan area is second only to Sydney's in population, and there is a good-natured rivalry between the two cities, to which geography and history have bequeathed diverse characteristics. Though Melbourne's flat site has led to the regular development of a rectangular pattern of streets, the city has many beautiful parks, and the person with an eye for architectural detail and history can find much that is varied and attractive. Melbourne has a reputation for conservatism and financial soundness—attributes that have contributed to its growth and are revealed by the burgeoning skyline of the central city and the rapidly expanding eastern suburbs.

This article is divided into the following sections:

Physical and human geography	829
The landscape	829
The city site	
Climate	
The city layout	
The people	829
Patterns of immigration	
Move to the suburbs	
The economy	830
Industry and trade	
Transportation	
Administration and social conditions	830
Government	
Health and education	
Cultural life	830
Arts	
Recreation	
History	831
Early settlement	831
Growth of the city	831
Bibliography	831

Physical and human geography

THE LANDSCAPE

The Melbourne Statistical Division, which includes all areas in close economic and social contact with the central city, covers 5,455 square miles (8,784 square kilometres). The continuous, densely settled metropolis, where a vast majority of the residents live, is less than half the size of the larger division, and the city proper is only 22.7 square miles (36.5 square kilometres).

The city site. Metropolitan Melbourne is situated at the northern end of Port Phillip Bay, 30 nautical miles (55 kilometres) from the bay's narrow entrance. Most of the flat terrain is less than 390 feet (120 metres) above sea level. The expansion of Melbourne from its origins at the mouth of the Yarra River to its present shape displays a strong correlation with the geology and drainage of the land. West of the original city site, basalt flows during the Tertiary Period filled the existing valleys and left flat, uniform plains. The eastern region, however, consists of undulating and dissected beds of sandstones, shales, and conglomerates laid down in the Silurian and Devonian periods. The thicker soils of the eastern region, together with its higher annual rainfall, supported a much denser cover of trees than on the basalt plains. Not surprisingly, the development of Melbourne has been mainly eastward into the broad reaches of land between Darebin Creek, the Plenty and Yarra rivers, and Koonung and Gardiners creeks. In a strikingly asymmetrical fashion, Melbourne's

urban development presently lines the entire eastern shore of Port Phillip Bay, from the mouth of the Yarra River to Point Nepean, 60 miles (97 kilometres) distant, while corresponding development on the west coast of the bay extends for only 10 miles.

Climate. Melbourne's weather results from the eastward flow of high-pressure cells separated by low-pressure troughs. These patterns follow a course that passes south of the continent in summer and over northern Victoria in winter. The annual rainfall of 26 inches (657 millimetres) is fairly evenly distributed throughout the year, with October usually the wettest month and January the driest. Temperatures are moderate, only rarely falling below freezing; average daily maximum temperatures vary from 55° F (13° C) in July to 79° F (26° C) in January. Winds associated with the eastward passage of weather systems ensure that Melbourne is spared the serious air pollution of some other large cities.

The city layout. The area of original settlement in Melbourne, which today forms its financial, legal, administrative, and ecclesiastical heart, was laid out in a rectangular pattern that has not changed. The area has a frontage along the Yarra River. Within this core are the major suburban and interstate railway stations, Victoria's Houses of Parliament, the Anglican and Roman Catholic cathedrals, arts and entertainment venues, museums, the Law Courts, the State Library, and many financial institutions, including the Melbourne Stock Exchange and the headquarters of major banks. Central to this area are two major destinations, Bourke and Swanston streets, which have been transformed into pedestrian malls, closed to automobile traffic. Most of the city's buildings are modern, but the Town Hall, the Law Courts, and the Exhibition Building provide excellent examples of 19th-century official architecture. The city is divided into 14 precincts, sectors identified by ethnic concentration, commercial clusters, or attractions.

The earliest suburbs—Carlton, Collingwood, Richmond, Prahran, St. Kilda, and Brighton—also have a generally rectangular pattern of streets. Row houses, often with verandas decorated with iron lacework, were common features in the suburbs close to the city's centre; and in parts of Carlton and South Melbourne some of these traditional terraces have been preserved.

THE PEOPLE

Patterns of immigration. The first official census of Melbourne, in 1836, numbered 177 persons, of whom 35 were

By courtesy of Promotion Australia



The Exhibition Building in Melbourne.

Asymmetrical urban growth

females. In the 1850s the gold rush in nearby areas of Victoria sparked the city's first major period of immigration. Newcomers came principally from other Australian colonies and Britain. By the 1920s Melbourne had become the home of more than half the residents of Victoria, and toward the end of World War II it reached a population of 1,000,000. This trend continued throughout the 20th century. By 2000 the Melbourne metropolitan area comprised nearly three-fourths of Victoria's population.

The second great wave of immigration came in the 1950s, when the Australian government pursued a deliberate policy of encouraging migration from Europe to provide workers for Australia's developing industries. The government provided assistance with travel costs and helped the immigrants to settle in Australia, learn English if necessary, and find employment. At first migrants were drawn mainly from the Baltic states and eastern Europe, many of these people being war refugees. Then larger numbers began coming from the United Kingdom and Ireland. Immigration agreements were signed with the Dutch, Maltese, West German, Italian, Greek, and Austrian governments. These programs sowed the seeds of Melbourne's present multicultural character. Italians, Greeks, and Yugoslavs formed the largest numbers of non-English-speaking migrants, and by the mid-1970s one-fifth of the city's population regularly spoke a language other than English. Immigration from Southeast Asia, in particular refugees from Vietnam and Cambodia, increased during the early 1980s.

Move to the suburbs. By 1966 the population of the Melbourne area had exceeded 2,000,000, and the city limits were rapidly being pushed eastward. The migrants from southern Europe were at first concentrated in the older industrial suburbs, where they found lower rents and nearby job opportunities. In those areas, too, support was available from earlier immigrants and the associations they had established. As these people prospered, however, they embraced the Australian dream of owning a home on a quarter-acre plot and began to move to the outer suburbs.

In the 1940s about 90 percent of Melbourne's residents lived within 10 miles of the central business district. By the 1980s this proportion had fallen to less than 50 percent, and the outer suburban areas had grown correspondingly. Although some western sections can be typified as working-class districts, on the whole Melbourne's suburbs represent a fairly even mixture of social classes and occupations. By the end of the 20th century, however, the trend began to reverse itself. The central city became increasingly residential as many aging buildings were rehabilitated and repurposed as housing. In addition, by 2000 the area's population had surpassed the 3,000,000 mark.

THE ECONOMY

Industry and trade. That Melbourne dominates the economic life of Victoria is not surprising—the city contains the bulk of the state's population. It is Victoria's financial centre and seat of government and is at the hub of the communications network linking the state to the rest of Australia and the world.

The city's original core offers the most employment, but employment is growing at a faster rate in the outer suburbs. The central city mainly holds service activities such as banking, insurance, retailing, entertainment, public accommodation, and railway transportation. Surrounding this core is an incomplete ring of inner industrial suburbs, where the first clothing and metal factories were established in the 19th century. In the outer suburbs, particularly to the east, small manufacturing areas began to develop after World War II, when these suburbs could offer large areas of inexpensive land, few problems of traffic congestion, and an increasing population.

In the 1990s an ambitious project was launched to develop Docklands, a 500-acre (200-hectare) site of crumbling industrial and port facilities, into a multiuse complex featuring high-technology businesses, parks and public spaces, restaurants, a theme park, and apartment buildings and other housing. Docklands was expected to become home for 15,000 people and a workplace for 20,000. The first

major facility, Colonial Stadium, a sports and entertainment venue, opened in 2000.

Melbourne's most important industries, in terms of numbers employed, are metal processing, including the manufacture of transportation equipment, and engineering. Other major industries include textile and clothing manufacture; food processing; papermaking and printing; and the manufacture of chemicals, furniture, and building materials. Melbourne is also one of Australia's leaders in the manufacture of computers and is developing as a centre for biomedicine and biotechnology.

The port of Melbourne occupies an area of level excavated land at the mouth of the Yarra River, southwest of the central business district. It is the nation's largest general-cargo port. The chief products handled are foodstuffs, crude oil and petroleum products, chemicals, and iron and steel.

Beginning in the 1960s, large regional shopping centres sprang up throughout Melbourne's outer suburbs, and the central city lost its dominant retail function. Nonetheless, major department stores and fashionable shops still exert a considerable pull on both residents and visitors.

Transportation. Melbourne is well served by an integrated public system of trains, buses, and tramcars, the latter a signature sight in the city. A network of national highways link Melbourne with adjoining states, and a system of freeways was greatly upgraded in the 1990s, including the creation of the Western Ring Road as a bypass route. The City Link project joined three major freeways with a bridge, tunnels, highway extensions, and interchanges to facilitate traffic movement. An underground rail loop serves the central business district. Melbourne's international and domestic airport is located at Tullamarine, 14 miles northwest of the city's centre.

ADMINISTRATION AND SOCIAL CONDITIONS

Government. The Victoria state government has the ultimate responsibility for Melbourne's major planning decisions and for providing its principal health, educational, and transport services. Local government in the Melbourne Statistical Division is provided by more than 30 entities. Councillors, led by a lord mayor, are elected under a system of compulsory voting conducted by mail. Councillors receive a stipend and represent the city at large. The councils pass local ordinances and control a number of services connected with building regulations, community welfare, garbage collection, and vehicle parking. Revenue for these purposes is raised by property taxes (rates).

Health and education. Since the system of free public hospitals was started in 1846, it has grown to encompass numerous special facilities, which deal with either particular ailments or categories of patients, as well as general hospitals. There are, in addition, many private hospitals.

The University of Melbourne, one of the oldest in Australia, was established in 1853 (though the first students were admitted in 1855); Monash and La Trobe universities were established in the 1960s, and Deakin University, established in 1974, maintains three campuses. Melbourne also has colleges of advanced education that offer degrees or diplomas in a variety of technical and academic subjects.

CULTURAL LIFE

Arts. Melbourne's already rich cultural life was greatly enhanced during the period between 1968 and 1984, when the Victorian Arts Centre was created on the south bank of the Yarra River close to the city's centre. It encompasses the National Gallery of Victoria, the Melbourne Concert Hall, and several theatres, among other facilities for the arts.

The National Gallery was originally opened in 1861 and moved to its present site in the arts centre in 1968. It houses several outstanding collections including, most notably, Australian art ranging from the colonial period to modern times; European art, with 18th-century works particularly well represented; and decorative arts.

The Melbourne Concert Hall seats 2,600. Its patrons appreciate not only the technical brilliance of the acoustic en-

Diverse ethnic population

The commercial port

The arts centre

gincering but also the hall's superb decorations in colours derived from the gemstone and mining industry, which makes the hall appear to have been carved out of a hillside.

Performance spaces of the arts centre include the State Theatre, the Playhouse, the George Fairfax Studio, and the Black Box, providing facilities for opera, ballet, musical theatre, drama, stand-up comedy, cabaret, and a variety of music. The Sidney Myer Music Bowl, in King's Domain Gardens opposite the arts centre, is an outdoor venue seating 13,000. Yearly seasons of opera, ballet, and concert music include performances by international artists. The Melbourne Symphony Orchestra was formed in 1949 and has toured North America, Japan, and New Zealand, as well as major Australian cities.

Recreation. Automobile license plates in Victoria carry the identification "Victoria—Garden State." Melbourne is worthy as the capital of a garden state, with more than one-fourth of its inner-city area consisting of public parks and reserves. These spaces were set aside in the mid-19th century, at a time when many civic leaders in other cities were concerned with commercial development rather than with the quality of life. Extensive tracts have also been allotted as parklands in the newer outer suburbs. The most famous park in Melbourne is the Royal Botanic Gardens (RBG). This area of 89 acres was established in 1846 and today contains lakes, lawns, and thousands of named trees and shrubs. The associated National Herbarium of Victoria, which houses a collection of some 1,200,000 pressed plant specimens, is internationally recognized and used by scholars. The RBG also maintains a separate 200-acre facility at Cranbourne, about 30 miles southeast of central Melbourne.

Melbourne has hundreds of sports fields, tennis courts, swimming pools, and golf courses for active sports participants. Spectators find good accommodations at the Melbourne Cricket Ground, which holds 100,000 and is used for both cricket and Australian Rules football, and at the Flemington Racecourse, where the valuable Melbourne Cup race is held every November. Melbourne hosted the 1956 Summer Olympic Games. New sports facilities, including a large tennis stadium near the Melbourne Cricket Ground, have been added since the Olympics. Sailing and fishing on Port Phillip Bay and surfing on the ocean beaches are also popular, and the winter ski slopes around Mount Buller are within easy reach.

History

EARLY SETTLEMENT

Port Phillip Bay was discovered by Europeans in 1802, when captains John Murray and Matthew Flinders visited the bay within a few months of each other. This area was then part of the colony of New South Wales, and the colony's governor, Philip Gidley King, instructed the surveyor-general, Charles Grimes, to examine the shores of the bay with a view to identifying sites for future settlement. In 1803 Grimes and his party discovered the Yarra River and traveled along its lower course. Unlike some members of the party, Grimes was not enthusiastic about the Yarra River as a potential settlement. Later in the same year Captain David Collins arrived with a contingent of soldiers and convicts and settled near Sorrento, just inside the entrance to the bay on the east coast. Within a few months, however, he decided that the location was unsuitable and moved his group to Tasmania.

Permanent settlement was delayed until 1835, when a pioneer settler and entrepreneur, John Batman, negotiated a treaty with the Aboriginal elders for the purchase of 500,000 acres at the head of Port Phillip Bay. The price was 40 blankets, 30 axes, 100 knives, 50 scissors, 30 mirrors, 200 handkerchiefs, 100 pounds of flour, and six shirts. Batman and his heirs were bound by the treaty to provide an annual "rent of tribute" of similar items.

A few days after the treaty was signed, Batman left, and two months later a party led by another pioneer, John Fawkner, settled on the banks of the Yarra River. There has been much debate about whether Batman or Fawkner should be regarded as the founder of Melbourne. Both seem to have an equal claim, but if the term is interpret-

ed to include expansion and consolidation of the settlement, then the honour must go to Fawkner. Within four years of signing his treaty, which had been disallowed by the Governor of New South Wales, Batman died, at age 38; his financial affairs were in disarray, and prolonged litigation over his will destroyed the estate he had created. Fawkner lived to the age of 76. He died in Melbourne in 1869 after a rewarding career in which he established hotels, a newspaper, and a bookselling business, acquired large areas of land, and held a seat on the Legislative Council for 18 years.

Melbourne is distinguished from the other Australian state capitals in that it was founded unofficially, by individual enterprise. Once Batman, Fawkner, and others had established the settlement in 1835, the government in Sydney had to recognize the fact. In 1836 the first administrator of the Port Phillip District arrived, and in 1837 the new settlement was given its present name honouring the British prime minister, William Lamb, 2nd Viscount Melbourne (of Kilmore). Melbourne became a town in 1842 and a city in 1847, but its first main surge in growth came in the early 1850s following the discovery of gold near Bendigo and Ballarat less than 100 miles away. In three years the population of Melbourne increased fourfold to 80,000.

The gold rush

GROWTH OF THE CITY

Melbourne capitalized on its central position within Victoria and its port facilities to capture most of the region's trade. Between 1856 and 1873, railways were built to Geelong, Ballarat, Bendigo, Echuca, and Wodonga, and in 1883 a link with the New South Wales rail system was established at Albury. In 1877 the Melbourne Harbour Trust was created, and the Coode Canal was cut in the soft alluvial sediments of the lower Yarra River to provide a more direct course free from silting problems.

During the 1870s manufacturing flourished under the protection of a high tariff, and progress in most spheres continued until 1889, when a financial crisis and the collapse of many firms lowered public confidence. The following decade witnessed a severe economic depression that began with a maritime strike and the failure of a number of banks and was sealed by seven years of drought from 1895 to 1902. In the decade before 1891 the population of Melbourne had increased by 200,000; in the following decade it rose by only 6,000.

In the early years of the 20th century, confidence gradually revived. Australia became a commonwealth, and Melbourne served as its federal capital until 1927, when Canberra was established. World Wars I and II encouraged the growth of manufacturing, and after 1945 European immigrants began to arrive in significant numbers.

After 1971 Melbourne's rate of growth slackened as immigration declined, and economic conditions worsened through the 1970s and early '80s. Nevertheless, during this period of slower population growth a number of major changes took place. The appearance of the inner city was transformed by the replacement of old buildings with multistory office structures and hotels. The system of arterial roads was improved dramatically. Several important suburban economic areas emerged, reducing the retail and industrial importance of the city's centre. And the cultural life of the city was immeasurably enlivened by the completion of the Victorian Arts Centre. The revitalization of central Melbourne continued after 1990 as the residential population grew and the massive Docklands development project began to transform the long-neglected waterfront area into a spectacular urban showplace.

BIBLIOGRAPHY. W.H. NEWNHAM, *Melbourne: Biography of a City* (1985), an excellent historical account; ANTHONY HARVEY (comp.), *The Melbourne Book* (1982), a good general description; J.S. DUNCAN (ed.), *Atlas of Victoria* (1982), which includes many well-illustrated references to Melbourne; and GRAEME DAVISON, *The Rise and Fall of Marvellous Melbourne* (1978), an economic history. See also AUSTRALIAN BUREAU OF STATISTICS, *Victorian Year Book* (annual), invaluable for statistics and developments in the area, with a good bibliography; and "Greater Melbourne," *International Demographics*, 5(3):1-9, 12 (March 1986), for recent demographic data on the city. (J.R.V.P./Ed.)

Memory

That experiences influence subsequent behaviour is evidence of an obvious but nevertheless remarkable activity called remembering. Learning could not occur without the function popularly named memory. Practice results in a cumulative effect on memory leading to skillful performance on the tuba, to recitation of a poem, and even to reading and understanding these words. So-called intelligent behaviour demands memory, remembering being prerequisite to reasoning. The ability to solve any problem or even to recognize that a problem exists depends on memory. Typically, the decision to cross a street is based on remembering numerous earlier experiences.

Practice (or review) tends to build and maintain memory for a task or for any learned material. Over a period of no practice what has been learned tends to be forgotten; and the adaptive consequences may not seem obvious. Yet, dramatic instances of sudden forgetting (as in amnesia) can be seen to be adaptive. In this sense, the ability to forget can be interpreted to have survived through a process of natural selection in animals. Indeed, when one's memory of an emotionally painful experience leads to severe anxiety, forgetting may produce relief. Nevertheless, an evolutionary interpretation might make it difficult to understand how the commonly gradual process of forgetting survived natural selection.

In speculating about the evolution of memory, it is helpful to consider what would happen if memories failed to fade. Forgetting clearly aids orientation in time; since old

memories weaken and the new tend to be vivid, clues are provided for inferring duration. Without forgetting, adaptive ability would suffer; for example, learned behaviour that might have been correct a decade ago may no longer be. Cases are recorded of people who (by ordinary standards) forgot so little that their everyday activities were full of confusion. Thus, forgetting seems to serve the survival of the individual and the species.

Another line of speculation posits a memory storage system of limited capacity that provides adaptive flexibility specifically through forgetting. In this view, continual adjustments are made between learning or memory storage (input) and forgetting (output). Indeed, there is evidence that the rate at which individuals forget is directly related to how much they have learned. Such data offer gross support of contemporary models of memory that assume an input-output balance.

Whatever its origins, forgetting has attracted considerable investigative attention. Much of this research has been aimed at discovering those factors that change the rate of forgetting. Efforts are made to study how information may be stored; that is, to discover the ways in which it may be encoded. Remembered experiences may be said to consist of encoded collections of interacting information; and interaction seems to be a prime factor in forgetting.

This article treats the phenomena and functions of normal memory; it also discusses the variety of disorders that affect the memory. The article is divided into the following sections:

The nature of memory: retention and forgetting	832
Measuring retention	832
Recall	
Recognition	
Relearning	
Time-dependent aspects of retention:	
storage and retrieval	833
Theories of forgetting	834
Correlates of rate of forgetting	835
Degree of learning	

Mnemonic systems	
Individual differences	
Abnormalities of memory	836
Organic disorders	836
Psychological studies of amnesia	838
Psychogenic amnesia	838
Paramnesia and confabulation	839
Hypermnesia	840
Bibliography	840

The nature of memory: retention and forgetting

Psychologists of the modern era, from their earliest speculations about remembering to the formulation of most of their latest experimentally based views, commonly have assumed that the critical problems are concerned with the physiological mechanisms by which events and experiences can be retained so that they can be mentally reproduced, either in their original mode or with the assistance of signs and symbols that are regarded as equivalent to that mode. Memory is thus usually considered to function perfectly in proportion to its literal accuracy of reduplication. Investigators have generally supposed that anything that influences the behaviour of an organism endowed with a central nervous system leaves—somewhere in that system—a "trace" or group of traces. So long as these traces last they can, in theory, be restimulated and the event or experience that established them will be remembered. The experimental psychology of remembering—all modern experts claim to base their conclusions upon experimental evidence—endeavours to discover methods for identifying the necessary and sufficient conditions for the persistence and length of persistence of traces and for their restimulation.

MEASURING RETENTION

Standard sentences, prose passages, and poems have been used to control input in studies of retention; but discrete

verbal units (such as words or sets of letters) are most frequently employed. The letters usually comprise lists of consonant syllables (three consonants; *e.g.*, RQK) or so-called nonsense syllables (consonant-vowel-consonant; *e.g.*, ROK). The order in which verbal units are to be learned and to be recited may be left to the subject (free recall). A schoolchild who can recite the names of all African countries probably has learned such a free-recall task. Units also can be presented serially (in a constant order), the subject being asked to recite them in that order; reciting the alphabet in the usual way represents such serial learning.

Pairs of words may be offered; in such paired-associate tasks the subject eventually is asked to produce the missing member of each pair when only one word is shown. This is akin to learning English equivalents for words from another language.

For these and similar tasks investigators commonly permit subjects enough practice trials to reach some preselected criterion or level of performance. This level effectively defines an immediate retention score against which later forgetting may be measured. Subsequent tests of retention are then made to investigate the rate at which forgetting proceeds. This rate tends to vary with the methods used, basically those of recall, recognition, or relearning.

Recall. Subjects may be asked to reproduce (recall) previously learned data in any order or in the original order in which they were learned.

Free recall,
serial
learning,
paired-
associate
tasks

In a free-recall test the instructions might be: "Yesterday you learned a list of words; please write as many of those words as you possibly can as they occur to you." For the paired-associate task the subject may be told: "Yesterday you learned some pairs of words; I will show you one word from each pair and you try to give the other." He may be paced, being limited to a few seconds to produce each word; or he may be unpaced, being given no rigidly specified limits.

If retention of any kind is to be measured over different periods (e.g., an hour, a day, a week) a separate group of individuals should be used for each period. The reason is that the very act of remembering constitutes practice that keeps memory lively, tending to give misleading underestimates of the rate of forgetting if the same subjects are tested over successive intervals.

Recognition. The subject's task is simpler in tests of recognition, since reproduction or retrieval (as in recall) is not required. The subject simply is asked to remember previously presented information when it is offered to him again. For example, he may be given a list of words for study; on the subsequent test of retention these are mingled with additional words, the subject being asked to identify (recognize) the original words. Apparently the recognition test stresses ability to choose between "old" (studied) data and "new" words, although this need not mean that choices are based only on temporal discrimination (awareness of time distinctions).

In an alternative variety of recognition test, each word studied might be paired with a new one, the task being to choose the old member of each pair. Or, the test words might be presented one at a time for identification as old or new. Sometimes learning and testing are combined: a very long list of words may be presented one at a time, some being repeated; the task is to recognize the repeats.

Some recognition tests stress memory of the order of presentation. The subject learns a serial list (reciting in a prescribed order); the list then is scrambled, and he is tested on his ability to rearrange it appropriately. Order may be based on how units are arranged in space (e.g., printed on a page) or on their numerical position in a series or on associative information. Thus, if a paired-associate list has been learned, the test may consist of the unmatched presentation of all units with a request to pair them properly. This sort of recognition seems to emphasize associative attributes. If some elements on the test were not presented originally, the temporal attribute also may be involved.

Relearning. The number of successive trials a subject takes to reach a specified level of proficiency may be compared with the number of trials he later needs to attain the same level. This yields a measure of retention by what is called the relearning method. The fewer trials needed to reach the original level of mastery, the better the subject seems to remember. The relearning measure sometimes is expressed as a so-called savings score. If 10 trials initially were required, and five relearning trials later produce the same level of proficiency, then five trials have been saved; the savings score is 50 percent (that is, 50 percent of the original 10 trials). The more forgetting, the lower the savings score.

Although it may seem paradoxical, relearning methods can yield both sensitive and insensitive measures of forgetting. Tasks have been devised that produce wide differences in recall but for which no differences in relearning are observed. (Some theorists attribute this to a form of heavy interference among learned data that has only momentary influence on retention.) Six months or a year after initial learning, some tests may give zero recall scores but can show savings in relearning.

When relatively long retention intervals (usually hours or days) are used, the methods are said to involve long-term memory. In a sense, methods for studying short-term memory are miniaturized versions of these. A list may be as short as one item, level of proficiency is very low, and retention intervals are in seconds (or minutes at most).

For example, the subject may be shown a single nonsense syllable for a few seconds' study. Next he is given a simple task (such as counting backward) to occupy him for a

half minute so that he cannot rehearse, and then is asked to recall the syllable. Forgetting is observed to occur over such short intervals, tending to be greater when length of interval increases, as in long-term memory. The same procedure can be used with a single paired-associate item or with a short list of four or five pairs. In a short-term counterpart of serial learning, a string of about eight single-digit numbers or letters is presented very rapidly (say, two per second), and the subjects are asked to recall them in the order in which they were presented. Recognition tests also can be adapted for measuring short-term retention. When only one presentation is used for learning, however, relearning measures are obviously unfeasible.

TIME-DEPENDENT ASPECTS OF RETENTION; STORAGE AND RETRIEVAL

Some workers theorize a distinct short-term memory system of sharply limited capacity that can retain information perhaps only a few seconds and a long-term system of relatively unlimited capacity and retention.

Among typical people, short-term function seems limited to about seven separate units (e.g., seven random letters or unrelated common words). Thus, one may consult a telephone directory and forget the number before dialing is completed. Information seems to enter long-term storage by such processing as rehearsal and encoding, as if short-term retention is a way station between incoming information and more enduring memory.

Other theorists do not distinguish short- and long-term systems as inferred from observed differences in capacity and retention. Positing only one storage system, they attribute short-term phenomena to very low levels of learning. Those who postulate distinct systems point to the results of injury to a specific brain region (the hippocampus): (1) information stored prior to hippocampal damage seems to be retained; (2) sufferers seem incapable of new long-term storage; (3) the short-term functions appear to be unimpaired and subjects perform as well as ever in tests of immediate memory (e.g., for a set of random numbers). It is as if new information no longer can be transferred from some sort of short-term system to relatively enduring storage.

Other data that bear on the controversy among theorists come from studies of people without known brain injury. When one has just seen a new list of words one at a time, the initial words in the list tend to be recalled best (primacy effects), those at the end next best (recency effects), while items from the middle are least likely to be recalled. This is quite consistently found as long as recall begins immediately following presentation of the last word. If, however, a short interval follows, during which the subject is otherwise occupied to prevent rehearsal, the recency effect may completely disappear: words at the end are no better recalled than those in the middle. Primacy effects are essentially undisturbed, while a delay as short as perhaps 15 seconds is enough to abolish the recency phenomenon. Although some suggest that recency effects depend on a separate short-term memory system and that primacy effects are mediated by a long-term system, a single memory function also may be invoked to accommodate the findings. Nevertheless, interest is growing in multisystem theories on the grounds that they enhance appreciation of the processes involved in establishing relatively enduring memory.

Investigators concerned with physiological bases for memory seek a kind of neurochemical code with enough stability physically to produce a structural change or memory trace (engram) in the nervous system; mechanisms for decoding and retrieval also are sought. Efforts at the strict behavioral level similarly are directed toward describing encoding, decoding, and retrieval mechanisms as well as the content of the stored information.

One way to characterize a memory (or memory trace) is to identify the information it encodes. A learner may encode far more information than is apparent in the task as presented. For example, if a subject is shown three words for a few seconds and (after 30 seconds of diversion from rehearsal) is asked to recall and then another triad of words is given under the same procedure, then

Neuro-
logical
evidence

Short-
term and
long-term
memory

another, and so on, then if all triads share some common element (e.g., all are animal names), poorer and poorer recall is observed on successive trials. Such findings may be explained by assuming that the learner encodes this animal category as part of his memory for each word. Initially, the common code might be expected to aid recall by sharply delimiting the word population. Successive triads, however, should tend to be encoded in increasingly similar ways, blurring their unique characteristics for the subject. An additional step provides critical supporting evidence for such an interpretation. If a final triad of vegetable names is unexpectedly presented, recall recovers dramatically. The subject tends to reproduce the vegetable names much better than he does those of the last animal triad; recall is about as efficient as it was for the first three animal names. This particular shift clearly seems to provide escape from earlier confusion or blurring, and it may be inferred that a common conceptual characteristic was encoded for each animal name.

Encoding
mech-
anisms

Any characteristic or attribute of a word may be investigated in this way to infer whether it is incorporated in memory. When recall does not recover it would seem that the manipulated characteristic has little or no representation in memory. For example, grammatical class typically does not appear to be encoded; decrement in recall produced after a series of triads consisting of verbs tends to continue when a shift is made to three adjectives. Such an experiment does not indicate what common encoding characteristic might be responsible for the decrement, suggesting only that it is not grammatical class.

Encoding mechanisms also may be inferred from tests of recognition. For example, subjects study a long list of words, being informed of a multiple-choice memory test to follow. Each word studied is made part of a test item that includes other carefully chosen new words (distractors). Distractors are selected to represent different types of encoding the investigator suspects may have occurred in learning. If the word presented for study is chosen by the subject, little can be inferred about the nature of the encoding. Any errors, however, can be most suggestive. Thus, if the word presented for study was TABLE, the multiple-choice item might be TABLE, CHAIR, ABLE, FURNITURE, PENCIL. If CHAIR is incorrectly selected, it may be suspected that this associatively related word occurred to the subject implicitly during learning and became so well encoded that the subject later could not determine whether it or TABLE had been presented. If the wrong choice is ABLE, acoustical resemblance to TABLE may have contributed to the confusion. If FURNITURE is erroneously chosen, perhaps conceptual category was prominent in the encoding.

Since it is not related in any obvious way to TABLE, the word PENCIL may be intended as a control, unlikely to be a part of the memory for TABLE. If this is the case, subjects should be more likely to select distractor words other than PENCIL (if indeed they have been encoded along with TABLE).

Although a subject may have encoded in ways suggested by particular distractors, he still may be able to choose the correct word. Or he may have encoded in ways not represented by the distractors.

Evidence has been accumulating to suggest that a long-term memory is a collection of information or of attributes that can serve in discriminating it from other memories and can function as retrieval cues. In addition to verbal attributes, visual images may be a part of the memory; emotional responses produced at the time the memory is established may be incorporated.

The common experience of having a name or word on the tip of the tongue seems related to specific perceptual attributes. In particular, people who report the "tip-of-the-tongue" feeling tend to identify the word's first letter and number of syllables with an accuracy that far exceeds mere guessing. There is evidence that memories may encode information about when they were established and about how often they have been experienced. Some seem to embrace spatial information; e.g., one remembers a particular news item to be on the lower right-hand side of the front page of a newspaper. Research indicates that the

"Tip-
of-the-
tongue"
phenom-
enon

rate of forgetting varies for different attributes. For example, memories in which auditory attributes seem dominant tend to be more rapidly forgotten than those with minimal acoustic characteristics.

If a designated (target) memory consists of a collection of attributes, its recall or retrieval should be enhanced by any cue that indicates one of the attributes. For example, on failing to recall the term horse (included in a list just seen), one may be told that there was an animal name among the words. Or he may be asked if an associate term (say, barn or zebra) helps him think of a word he missed. While some additional recall has been observed with this kind of help, failures are common even with ostensibly relevant cues. Though it is possible that the cues frequently are inappropriate, nevertheless, if words were not learned (encoded or stored) with accompanying attributes, cuing of any kind should be ineffective.

THEORIES OF FORGETTING

When memory of past experience is not activated for days or months, forgetting tends to occur; and any theory of forgetting must cope with this primitive observation. Such auxiliary phenomena as differences in the rates of forgetting for different kinds of information also must be accommodated.

It has been theorized that as time passes the physiological bases of memory tend to change. With disuse, it is held that the neural engram (the memory trace in the brain) gradually decays or loses its clarity. While such a theory seems reasonable, it would, if left at this point, do little more than restate behavioral evidence of forgetting at the nervous-system level. Decay or deterioration does not seem attributable merely to the passage of time; some underlying physical process needs to be demonstrated. Until a neurochemical basis for memory can be more explicitly described, any decay theory of forgetting must await detailed development.

A pre-eminent theory of forgetting at the behavioral level is anchored in the phenomena of interference; in what are called retroactive and proactive inhibition. In retroactive inhibition, new learning interferes with retention of the old; in proactive inhibition, old memories interfere with the retention of new ones. Both phenomena have great generality in studies of any kind of learning, although most research among humans has considered verbal learning.

People may, for example, learn two successive verbal lists; the next day some are asked to recall the first list and others to recall the second. Still a third (control) group learns only one list and is asked to recall it a day later. People who learn two lists almost unfailingly will recall much less than do those in the control group. The amount by which controls exceed those who recall the first list is a measure of retroactive inhibition; the degree to which they are better than those who recall the second list is a measure of proactive inhibition. While retroactive inhibition usually will be observed in relearning, it is unusual to detect proactive deficit under such circumstances.

Theorists attribute the loss produced by these procedures to interference between list-learning tasks. When lists are constructed to exhibit varying differences, the degree of interference seems to be related to the amount of similarity. Thus loss in recall will be reduced when two successive lists have no identical terms. Maximum loss generally will occur when there appears to be heavy (but not complete) overlap in the memory attributes for the two lists. One may recall parts of the first list in trying to remember the second and vice versa. (This breakdown in discrimination may reflect the presence of dominant attributes that are appropriate for items in both lists.) Discrimination tends to deteriorate as the number of lists increases, retroactive and proactive inhibition increasing correspondingly, suggesting interference at the time of recall.

In retroactive inhibition, however, all of the loss need not be attributed to competition at the moment of recall. Some of the first list may be lost to memory in learning the second; this is called unlearning. If one is asked to recall from both lists combined, first-list items are less likely to be remembered than if the second list had not been learned. Learning the second list seems to act backward in

Inter-
ference
theory;
retroactive
and pro-
active
inhibition

time (retroactively) to destroy some memory for the first. So much effort has been devoted to studying conditions that affect unlearning that it has become a major topic in interference theory.

Retroactive and proactive effects can be quite gross quantitatively. If one learns a list one day and tries to recall it the next, learns a second list and attempts recall for it the following day, learns a third and so on, recall for each successive list tends to decline. Roughly 80 percent recall may be anticipated for the first list; this declines steeply to about 20 percent for the tenth list. Learning the earlier lists seems to act forward in time (proactively) to inhibit retention of later lists. These proactive phenomena indicate that the more one learns the more rapidly one will forget. Similar effects can be demonstrated for retroactive inhibition within just one laboratory session.

Such powerful effects have led some to theorize that all forgetting is produced by interference. Any given memory is said to be subject to interference from others established earlier or subsequently. Interference, theoretically, may occur when memories conflict through any attributes. With a limited group of attributes and an enormous number of memories, it might seem that everyday attempts to recall would be chaotic. Yet even if all of the memories shared some information, other attributes not held in common could still serve to distinguish them. For example, every memory theoretically is encoded at a different time and temporal attributes might serve to discriminate otherwise conflicting memories. Indeed, when two apparently conflicting lists are learned several days apart, proactive inhibition is markedly reduced. Assuming memories to be multiple-encoded, interference theory need not predict utter confusion in remembering.

Sources of interference are most pervasive and should not be considered narrowly. For example, any memory seems to be established in specific surroundings or context, and subsequent efforts to remember tend to be less effective when the circumstances differ from the original. Alcoholics, when sober, tend to have trouble finding bottles they have hidden while intoxicated; when they drink again, the task is much easier. Some contexts also may be associated with other memories that interfere with whatever it is that one is trying to remember.

Each new memory tends to amalgamate information already in long-term storage. Encoding mechanisms invariably adapt or relate fresh data to information already present, to the point that what is coded may not be a direct representation of incoming stimuli. This is particularly apparent when input is relatively meaningless; the newly encoded memory comes to resemble those previously established (*i.e.*, it accrues meaning). For example, a nonsense word such as LAJOR might be encoded as MAJOR.

To recall any nonsense word correctly requires that an appropriate decoding rule be a part of the memory; but coding rules are subject to forgetting (interference) in the same way that any attribute is. Qualitative changes in memory may result when the information presented does not allow precise decoding; thus, when one sees a drawing of a jagged figure that resembles a star he might encode it as a star, knowing full well that it is not perfect. Subsequent decoding in recall (or recognition) thus produces only an approximation of the original jagged figure; it may well be influenced by other equally imprecise decoding rules already stored. In like fashion, somewhat incoherent sentences may become more reasonable during encoding; they tend to be reproduced in memory tests more coherently. When the learner has trouble making sense of any new stimulus (when he cannot specify encoding and decoding rules with precision) the decoded memory tends to resemble previously established memories.

Although interference has attracted wide support as an account of forgetting, it must be placed in perspective. Interpretations that emphasize distinctions between short- and long-term memory and that posit control processes for handling information are potentially more comprehensive than is interference theory. Behavioral evidence for interference eventually may be explained within such systems.

In addition, a number of deductions from interference

theory have not been well supported by experiment. The focus of difficulty lies in the hypothesis that interference from established memories is a major source of proactive inhibition. The laboratory subject is asked to learn tasks with attributes that have varying degrees of conflict with memories established in daily life. Theoretically, the more conflict, the greater the proactive interference to produce forgetting. Yet a number of experiments have failed to provide much support for this prediction.

Interference theory also fails to account for some pathological forms of forgetting. Repression as observed in psychiatric practice, for example, represents almost complete, highly selective forgetting, far beyond that anticipated by interference theorists. Attempts to study repression through laboratory procedures have failed to yield systematic data that could be used to test theoretical conclusions.

CORRELATES OF RATE OF FORGETTING

Although forgetting normally is expected to begin as soon as practice ceases, at times an exception (known as reminiscence) has been reported. In reminiscence, memory seems to improve without practice; retention is even better if tested after a rest period than if tested immediately after learning trials stop. Observed only over periods of a few minutes, this elusive phenomenon produces very small improvements, and forgetting follows. Scores of studies designed to elicit reminiscence have failed to do so, yielding only evidence of forgetting. Reasons for the conflicting findings have not been identified.

Degree of learning. The degree of learning is found to be directly associated with the amount of practice. In a metaphoric sense, specific memory may be said to grow stronger and stronger as practice proceeds. Even after a task can be performed or recited perfectly, continued practice (sometimes called overlearning) increases the "strength" of the memory. The rate of forgetting is slower when the degree of learning is greater. If there were one universal prescription for resisting forgetting, it would be to learn to a very high level initially; results seem even better when learning trials are not bunched together. Practice trials may be given en masse in a single session or the same number of trials may be distributed in sessions held on different days. The interrupted schedule is far superior to massed practice in that the rate of forgetting that follows distributed practice is much slower. The laboratory evidence also confirms the belief that cramming for an examination may produce acceptable performance shortly afterwards, but that such massed study results in poor long-term retention. Information learned in widely distributed practice appears less susceptible to interference; memories established under distributed schedules also are less likely to produce proactive inhibition than are those learned in massed trials.

Mnemonic systems. The principle that new information is encoded to previously stored data has been used in an effort to aid memory function. When encoding techniques are formally applied, they are called mnemonic systems or devices. (The popular rhyme that begins "Thirty days hath September . . ." is an example.) Verbal learning can be enhanced by providing an appropriate mnemonic system (even to a bright college student who may have devised efficient systems of his own). Thus, paired associates (*e.g.*, DOG-CHAIR) will be learned more rapidly if they are included in a simple sentence (*e.g.*, The dog jumped over the chair). Imagery that can relate different words to be learned (even in a bizarre fashion) has been found beneficial. Some investigators hold that pure rote learning (in which no use is made of established memories except to directly perceive the stimuli) is rare or nonexistent. They suggest that all learning elaborates on memories already available. This could be taken to mean that the rate of forgetting would be the same whether or not a formal mnemonic system were used in learning.

Indeed, there seem to be no experimental results in which formal mnemonic instruction has resulted in forgetting more rapidly than when such special training is not given. Yet, while there often is no difference in the rate of forgetting, a number of studies indicate slower forgetting following instruction in a mnemonic system. These

Limits
on inter-
ference

Remi-
niscence

discrepancies may mean that some mnemonic systems are more subject to interference than are others. Perhaps the methods used fail to adequately distinguish between learning and forgetting.

Rates
of
learning
and
forgetting

Factors that influence the rate of learning should be distinguished from those that affect the rate of forgetting. For example, nonsense syllables are learned more slowly than are an equal number of common words; if both are studied for the same length of time, the better learned common words will be forgotten more slowly. But this does not mean that rate of forgetting *intrinsically* differs for the two tasks. Degree of learning must be held constant before it may be judged whether there are differences in rate of forgetting. Rates of forgetting can be compared only if tasks are learned to an equivalent degree. Indeed, when degree of learning is experimentally controlled, different kinds of information are forgotten at about the same rate. Nonsense syllables are *not* forgotten more rapidly than are ordinary words. In general, factors that seem to produce wide differences in rate of learning show little (if any) effect on rate of forgetting. (Despite discrepant evidence, mnemonic systems may prove an exception.)

Individual differences. Experimental findings seem to contradict the common intuition that people inherently differ in the rate at which they forget. This intuitive belief appears largely to derive from definite, wide individual differences in rate of learning; some people do learn faster than others. Thus, given the same number of trials or identical time in which to study, people will vary widely in the level of learning they achieve. Individual differences in forgetting then can be predicted efficiently, merely on the basis of how well each person has learned. This powerfully indicates that ordinary estimates of one's rate of forgetting are spurious, being obscured by uncontrolled differences in learning ability. One's talent for learning seems to swamp efforts to assess his inherent tendency to forget. Under less ordinary circumstances, however (*e.g.*, selective brain injury, stroke, neurotic amnesia), the degree of learning does seem to be almost completely irrelevant to the rate at which one forgets. An amnesia sufferer may forget his own name and still may be able to remember that Sofia is the capital of Bulgaria (see below). (B.J.U./Ed.)

Abnormalities of memory

Disorders of memory must have been known to the ancients and are mentioned in several early medical texts; but it was not until the closing decades of the 19th century that serious attempts were made to analyze them or to seek their explanation in terms of brain disturbances. Of the early attempts, the most influential was that of a French psychologist, Théodule-Armand Ribot, who, in his *Diseases of Memory* (1881, English translation 1882), endeavoured to account for memory loss as a symptom of progressive brain disease by embracing principles describing the evolution of memory function in the individual, as offered by an English neurologist, John Hughlings Jackson. Ribot wrote:

Ribot's
"law"

The progressive destruction of memory follows a logical order—a law. It advances progressively from the unstable to the stable. It begins with the most recent recollections, which, being lightly impressed upon the nervous elements, rarely repeated and consequently having no permanent associations, represent organization in its feeblest form. It ends with the sensorial, instinctive memory, which, having become a permanent and integral part of the organism, represents organization in its most highly developed stage.

The statement, amounting to Ribot's "law" of regression (or progressive destruction) of memory, enjoyed a considerable vogue and is not without contemporary influence. The notion has been applied with some success to phenomena as diverse as the breakdown of memory for language in a disorder called aphasia and the gradual return of memory after brain concussion. It also helped to strengthen the belief that the neural basis of memory undergoes progressive strengthening or consolidation as a function of time. Yet students of retrograde amnesia (loss of memory for relatively old events) agree that Ribot's principle admits of many exceptions. In recovery from

concussion of the brain, for example, the most recent memories are not always the first to return. It has proved difficult, moreover, to disentangle the effects of passage of time from those of rehearsal or repetition on memory.

A Russian psychiatrist, Sergey Sergeevich Korsakov (Korsakoff), may have been the first to recognize that amnesia need not necessarily be associated with dementia (or loss of the ability to reason), as Ribot and many others had supposed. Korsakov described severe but relatively specific amnesia for recent and current events among alcoholics who showed no obvious evidence of shortcomings in intelligence and judgment. This disturbance, now called the Korsakoff syndrome, has been reported for a variety of brain disorders aside from alcoholism and appears to result from damage in a relatively localized part of the brain.

The neurological approach may be combined with evidence of psychopathology to enrich understanding of memory function. Thus, a French neurologist, Pierre Janet, described amnesia sufferers who were apparently very similar to those observed by Korsakov but who gave no evidence of underlying brain disease. Janet also studied people who had lost memory of extensive periods in the past, also without evidence of organic disorder. He was led to regard these amnesias as hysterical, explaining them in terms of dissociation: a selective loss of access to specific memory data that seem to hold some degree of emotional significance. In his experience, reconnection of dissociated memories could as a rule be brought about by suggestion while the sufferer was under hypnosis. Freud regarded hysterical amnesia as arising from a protective activity or defense mechanism against unpleasant recollections; he came to call this sort of forgetting repression, and he later invoked it to account for the typical inability of adults to recollect their earliest years (infantile amnesia). He held that all forms of psychogenic (not demonstrably organic) amnesia eventually could resolve after prolonged sessions of talking (psychotherapy) and that hypnosis was neither essential nor necessarily in the amnesiac's best interest. Nevertheless, hypnosis (sometimes induced with the aid of drugs) has been widely used in the treatment of hysterical amnesia, particularly in time of war when only limited time is available.

Repression

ORGANIC DISORDERS

Defect of memory is one of the most frequently observed symptoms of impaired brain function. It may be transitory, as after an alcoholic bout or an epileptic seizure; or it may be enduring, as after severe head injury or in association with brain disease. When there is impaired ability to store memories of new experiences (up to total loss of memory for recent events) the defect is termed anterograde amnesia. Retrograde loss may progressively abate or shrink if recovery begins, or it may gradually enlarge in scope, as in cases of progressive brain disease. Minor grades of memory defect are not uncommon aftereffects of severe head injury or infections such as encephalitis; typically they are shown in forgetfulness about recent events, in slow and insecure learning of new skills, and sometimes in a degree of persistent amnesia for events preceding the illness.

Transient global amnesia. Apparently first described in 1964, transient global amnesia consists of an abrupt loss of memory lasting from a few seconds to a few hours, without loss of consciousness or other evidence of impairment. The individual is virtually unable to store new experience, suffering permanent absence of memory for the period of the attack. There is also a retrograde loss that may initially extend up to years preceding the attack. This deficit shrinks rapidly in the course of recovery but leaves an enduring gap in memory that seldom exceeds the three-quarters of an hour before onset. Thus the person is left with a persisting memory gap only for what happened during the attack itself and in a short period immediately preceding. Such attacks may be recurrent, are thought to result from transient reduction in blood supply in specific brain regions, and sometimes presage a stroke.

Traumatic amnesia. On recovery of consciousness after trauma, a person who has been knocked out by a blow on the head at first typically is dazed, confused, and imper-

fectly aware of his whereabouts and circumstances. This so-called posttraumatic confusional state may last for an hour or so up to several days or even weeks. While in this condition, the individual appears unable to store new memories; on recovery he commonly reports total amnesia for the period of altered consciousness (posttraumatic amnesia). He also is apt to show retrograde amnesia that may extend over brief or quite long periods into the past, the duration seeming to depend on such factors as severity of injury and the sufferer's age. In the gradual course of recovery, memories are often reported to return in strict chronological sequence from the most remote to the most recent, as in Ribot's law. Yet this is by no means always the case; memories seem often to return haphazardly and to become gradually interrelated in the appropriate time sequence. The amnesia that remains seldom involves more than the events that occurred shortly before the accident though in severe cases careful inquiry may reveal some residual memory defect for experiences dating from as long as a year before the trauma. It is thought by some that, after recovery, the overall period of time for which there is no recollection may indicate the degree of severity of the head injury.

Traumatic automatism. Posttraumatic amnesia is sometimes observed after mild head injury without loss of consciousness and with no apparent change in ordinary behaviour. A football player who is dazed but not knocked out by a blow on the head, for example, may continue to play and even score a goal. But he may be going through these motions automatically and may later have no memory whatever of the part of the game that followed his injury. The phenomenon is known as traumatic automatism and seems similar to, if not identical with, transient global amnesia.

Memory defect after electroconvulsive therapy. Electroconvulsive treatments have been widely used in psychiatry, particularly for depressed people. A seizure or convulsion is induced by passing current through electrodes placed on the forehead. Each treatment is followed by a period of confusion for which the person is subsequently amnesic; at this time there is also a rapidly abating amnesia of some seconds for events that immediately preceded the shock. After a number of treatments, however, some individuals complain of more persistent memory defect, shown mainly in exaggerated forgetfulness for day-to-day events. These difficulties nearly always clear up within a few weeks after treatment ends. Experimental evidence tentatively suggests that electroshock administered to only one side of the head produces therapeutic results equal to those of the standard procedure but with significantly reduced impairment of memory.

Korsakoff's syndrome. First described in cases of chronic alcoholism, Korsakoff's psychosis, or syndrome, occurs in a wide variety of toxic and infectious brain illnesses, as well as in association with such nutritional disorders as deficiency of the B vitamins. The syndrome also has been observed among people with cerebral tumours, especially those involving the third ventricle (one of the fluid-filled cavities in the brain). The main psychological feature is gross defect in recent memory, sometimes so severe as to produce "moment-to-moment" consciousness; such people can store new information only for a few seconds and report no continuity between one experience and the next. They seem incapable of learning, even after many trials or repetitions. Although cases of such severity are relatively rare, the ability to store experience only briefly is quite characteristic of Korsakoff's syndrome.

In addition, sufferers almost always show evidence of retrograde amnesia that can span as little as a few weeks past to as much as 15 or 20 years before onset of the disorder. These extensive retrograde amnesias are seldom total or uniform, and "islands" of memory often can be found by persistent interrogation. The person's memory function depends heavily on circumstances; for example, a man with Korsakoff's syndrome who recognizes his wife instantly when she visits may in her absence vehemently deny that he is married. Commonly, there is disorientation in place and time; the individual often underestimates his own age, sometimes grossly. Some sufferers characteris-

tically confabulate; *i.e.*, they remember experiences they never personally had or they falsely localize their memories in time. Sufferers sometimes deny their illness or memory problems. Otherwise, they can exhibit good intelligence and, apart perhaps from some lack of spontaneity, may show little or no personality change.

While Korsakoff's syndrome is commonly encountered as a transitory sign of brain disorder, it can be chronic, remaining effectively unimproved over many years. Even with improvement, however, an appreciable weakness in recent memory, particularly in regard to sequence in time, is quite apparent.

Persistent defect after encephalitis. Attention repeatedly has been drawn to severe and persistent memory defect following attacks of a form of brain inflammation called acute inclusion body encephalitis. The individual's behaviour closely resembles that of Korsakoff's syndrome except that his insight into the memory disorder is usually good and confabulation is infrequent or absent. Indeed, the memory disorder is sometimes so limited and specific as to raise the possibility of a psychogenic (*i.e.*, hysterical) amnesia. In cases of this kind there may be little or no impairment of intelligence or judgment.

Defect following brain surgery. Surgical operations on the sides of the brain (the temporal lobes) to remove tissues that produce symptoms of epilepsy are routine. While good results are often achieved, a degree of memory defect ensues. Operations on the dominant (usually left) temporal lobe tend to hamper one's ability to learn verbal information by hearing or reading. Usually observable even before surgery, the defect tends to be more marked after operation and has been reported to persist for up to three years before eventual recovery. Operations on one temporal lobe when there is unsuspected damage to its fellow on the other side of the brain (or on both lobes, in surgery very rarely undertaken) produce severe and persistent general memory defect, altogether comparable to postencephalitic amnesia. There is gross defect in recent memory and in learning (except perhaps in motor learning), with retrograde amnesia that initially may involve several years of the person's past. Intelligence otherwise appears to be well preserved; the individual shows insight into his memory difficulty, and seldom, if ever, confabulates.

Diffuse brain diseases. Some memory failure is almost universal during old age, particularly in forgetfulness for names and in the reduced ability to learn. Many people of advanced age, nevertheless, show adequate memory function if they suffer no brain disease. Impairment of memory is a characteristic early sign of senility, as well as of hardening of the brain arteries (cerebral arteriosclerosis) at any age, with exaggerated forgetfulness for recent events and progressive failure in memory for experiences that preceded the disorder. As arteriosclerotic brain disease progresses, amnesia tends to extend further into the past, embracing personal experience and general or common information. When the symptoms are almost those of Korsakoff's syndrome, the disturbance is called presbyphrenia. In most cases the amnesia is complicated by failure in judgment and changes in character. It has been suggested that severe memory defect in an elderly person carries a poor prognosis, being related to such factors as a shortened survival time and an increased death rate.

A Swiss psychiatrist, Eugen Bleuler, held that amnesia results only from a diffuse disorder of the outer layers (cortex) of the brain and suggested that memory depends on the integrity of the cortex as a whole. Indeed, the removal of brain tissue from rats and monkeys in experimental studies has indicated that retention of complex habits by the animals depends on the total amount of cortex that remains. It was claimed that the degree to which memory is lost depends not on where the brain is injured but on the extent of the damage. (This is the "law" of mass action, which asserts that the brain functions in a unitary manner; *i.e.*, as a whole.) While the extent of diffuse brain damage is roughly related to the severity of memory defect, the principle of mass action is manifestly inadequate. Whatever its physical basis, memory seems to depend on the integrity of relatively limited parts of the brain, rather than on that organ (or even the cortex) as a whole.

Haphazard recovery of memory

Extensive retrograde amnesia

Senility

Brain structures in amnesia

Severe and highly specific amnesic symptoms principally stem from damage to such brain structures as the mammillary bodies, circumscribed parts of the thalamus, and of the temporal lobe (e.g., the hippocampus). While the ability to store new experience (and perhaps to retrieve well-established memories) appears to depend on a distinct neural system involving the temporal cortex and limited parts of the thalamus and hypothalamus, understanding of the neuroanatomy of memory remains sketchy enough to generate major differences of opinion. French and German workers tend to stress the role of the mammillary bodies, while U.S. investigators tend to implicate the thalamus. It has been pointed out that circumscribed damage to the mammillary bodies is not invariably associated with memory defect; cases of amnesia evidently occur in which these structures are spared. Nevertheless, implication of the mammillary bodies in a large number of verified cases of Korsakoff's syndrome seems incontrovertible. Injury to other neural tissues (e.g., the so-called fornix bundle deep within the brain) that anatomically might be expected to produce severe memory disorder rarely does so. While evidence for amnesia as a sign of localized brain damage is impressive, much remains to be understood about the physical system that sustains memory.

PSYCHOLOGICAL STUDIES OF AMNESIA

Short-term memory. The so-called short-term memory is typically intact among amnesia sufferers. Such victims usually can repeat a short phrase or a series of words or numbers from immediate memory as adequately as anyone of comparable age and intelligence. Such an amnesic person can retain the gist of a question or request long enough to respond appropriately, unless, of course, there is enough delay in performance or attention is diverted. Evidently the ability to register information is intact, if this means availability of data in short-term memory. Thus, experimental psychologists who favour a sharp distinction between short-term and long-term storage systems contend that the primary deficit in amnesia is an inability to transfer information from short-term to long-term storage.

Associative learning. It has been argued that the basic deficit in the amnesic state is a loss of learning ability. In a series of experiments with amnesic patients, using, for the most part, verbal material, the subjects evidenced failure to link new with old associations, rapid fading of new associations, and great difficulty in reproducing whatever associations might have been formed. These findings have been amply confirmed. In one view, however, the weakness resides less in the failure to establish new associations than in their rapid decay (i.e., accelerated forgetting). On the other hand, it has often been noticed that if a Korsakoff patient can once succeed in learning an item, he may be able to reproduce it correctly after an appreciable interval of time. Further experiments, using a variety of techniques for assessing learning and retention, have suggested that retrieval rather than learning is at fault.

Motor skill. It has been noted that the fact that the acquisition of manual skill in Korsakoff patients is less impaired than either verbal learning or the solution of puzzles or mazes. This is confirmed in the observation that a severely amnesic patient who had undergone an extensive operation on the temporal lobes could perform rotary-pursuit and tracking tasks at a level not greatly inferior to that of healthy subjects. A second case of the same kind has been described, in which memory for motor tasks such as maze learning or the rendering of new compositions on the piano is said to have been completely preserved. These observations suggest that the acquisition of motor skill may remain relatively unaffected by lesions that give rise to a severe defect of general memory. What is commonly called global or generalized memory defect may, therefore, become increasingly subject to fractionation.

Residual learning capacity. Korsakov himself pointed out that a patient who consistently denies having seen his doctor before does not necessarily react to him on each successive encounter as a total stranger. It thus appears that, despite gross amnesia, some learning, perhaps implicit, can still take place. This view has gained much support from clinical and experimental studies. About 1900 it

was reported that even severely affected Korsakoff patients show appreciable savings in relearning verbal material after an interval of several hours or days, thus indicating minimal retention. Some Korsakoff patients, in spite of gross amnesia, eventually learn their way about the hospital. Again, some patients who disown any knowledge of their whereabouts may nevertheless give the correct name of the hospital, when asked to guess or to select it from a list containing the names of several hospitals. Thus, while learning capacity is seldom, if ever, wholly destroyed, there is failure to integrate new knowledge within the total personality. It is apparently a lack of mental cohesion that lies at the basis of Korsakoff's psychosis.

Forgetting. While some clinicians have attributed memory defect largely to defective registration of experience (i.e., failure to form memory traces), the widely accepted view is that it results primarily from a greatly increased rapidity of forgetting (i.e., rapid decay of memory traces). This view has also been held by the great majority of experimental psychologists who have worked with amnesic people. The consensus is that amnesia suffers characteristically lose much of the memory they once had. This conclusion finds support in the very rapid extinction of conditioned eyeblink responses to a buzzer. It is notable that, in Korsakoff states, forgetting appears to be due to the passage of time (oblivescence) rather than to retroactive inhibition or some kindred interference effect.

Time disorders. Estimation of time is typically poor in amnesic states. The individual is prone to underestimate grossly the time in which he has been engaged on any particular activity. Conversely, he may equally grossly overestimate the time that has elapsed since a particular event (e.g., the visit of a relative) of which he has preserved some recollection. Indeed, amnesic patients exhibit a remarkable want of coherence in their thought processes, suggesting that a lack of temporal synthesis underlies, and may indeed in large part explain, the defect of memory. Yet although difficulties in dating particular past events and in building a coherent time framework are characteristic of amnesic states and may persist after otherwise good recovery, an explanation couched wholly in terms of time disturbance is scarcely convincing.

Retrograde amnesia. Since retrograde amnesia relates to memory for events that took place when brain function was unimpaired, it clearly cannot be ascribed to failure of registration—with the exception, perhaps, of the very brief permanent amnesias following electroconvulsive shock or head injury. Retrograde amnesia otherwise would appear to be wholly due to a failure of retrieval, though this failure is evidently selective. That recent memories are generally harder to evoke than those more remote is usually explained on the basis of consolidation; i.e., progressive strengthening of memory traces with the passage of time. Yet, recency is not the only factor, and in some cases memory for a relatively recent event may still be preserved while that for one more remote is inaccessible. Much depends, too, on the method used to test retrieval; e.g., recognition may succeed when voluntary recall entirely fails. By and large, the availability of information in memory would seem to depend to a considerable extent on its relation to the person's current interests and preoccupations. When these are severely curtailed by an amnesic state, the links connecting present and past are severed, with a consequent failure of reproduction.

PSYCHOGENIC AMNESIA

Some forms of amnesia appear to be quite different from those associated with detectable injury or disease of the brain. These comprise, first, amnesias that can be induced in apparently normal individuals by means of suggestion under hypnosis; and secondly, amnesias that arise spontaneously in reaction to acute conflict or stress, and which are commonly called hysterical. Such amnesias are reversible and have been explained wholly in psychological terms. Nevertheless, organic factors are not infrequently involved to some extent, and the distinction between organic and psychogenic amnesia may turn out to be far less absolute than has been supposed.

Hypnotic amnesia. Memory of a hypnotic trance is

Generalized versus specific memory loss

Consolidation of the memory trace

often vague and fragmentary, as in awakening from an ordinary dream. This may be due in part to defect of registration during the period of altered consciousness. At the same time, very much more complete posthypnotic amnesia can be induced if an individual is told that, when he awakens, he will remember nothing of what went on during the period of hypnosis. This is clearly a psychogenic phenomenon; memory is fully regained if the patient is rehypnotized and an appropriate counter-suggestion given. It may also be regained if the person is persistently interrogated in the waking state, again suggesting that the amnesia is apparent rather than real. This observation led Freud to seek access to ostensibly forgotten (repressed) memories in his patients without the use of hypnosis.

Hysterical amnesia. Hysterical amnesia is of two main types. One involves the failure to recall particular past events or those falling within a particular period of the patient's life. This is essentially retrograde amnesia but it does not appear to depend upon an actual brain disorder, past or present. In the second type there is failure to register—and, accordingly, later to recollect—current events in the patient's ongoing life. This is essentially anterograde amnesia and, as an ostensibly psychogenic phenomenon, would appear to be rather rare and almost always encountered in cases in which there has been a preexisting amnesia of organic origin. Rarely, amnesia appears to cover the patient's entire life, extending even to his own identity and all particulars of his whereabouts and circumstances. Although most dramatic, such cases are extremely rare and seldom wholly convincing. They usually clear up with relative rapidity, with or without psychotherapy.

Hysterical amnesia differs from organic amnesia in important respects. As a rule it is sharply bounded, relating only to particular memories, or groups of memories, often of direct or indirect emotional significance. It is also usually motivated in that it can be understood in terms of the patient's needs or conflicts; *e.g.*, the need to seek financial compensation after a road accident causing a mild head injury or to escape the memory of an exceptionally distressing or frightening event. Hysterical amnesia also may extend to basic school knowledge, such as spelling or arithmetic, which is never seen in organic amnesia unless there is concomitant aphasia or a very advanced state of dementia. A most distinctive feature of hysterical amnesia is that it can almost always be relieved by such procedures as hypnosis. Although distinguishing organic from psychogenic amnesia is not always easy, it can usually be achieved on the basis of such criteria, especially when there is no reason to suspect actual brain damage.

Legal implications. The differentiation of organic from functional amnesia not uncommonly assumes legal importance, as in cases in which compensation is sought for disability held to be due to industrial or road accidents causing head injuries. If there is a complaint of defective memory, it is legally important to ascertain what part of it can be ascribed to the aftereffects of the head injury and what part of it to subsequent psychogenic elaboration. Similar issues may also arise on occasion in criminal cases, as in a trial in England (1959) in which it was contended that the accused man had a total amnesia for the circumstances of his alleged offense—the murder of a police officer—and should therefore be regarded as unfit to plead. After much discussion as to whether the amnesia was organic, hysterical, or feigned, the jury found it not to be genuine and the trial proceeded to conviction.

Mixed amnesic states. Students of amnesia have been increasingly impressed by the frequency with which psychogenic factors appear to reinforce, prolong, or otherwise complicate an organic memory defect. Hysterical reactions appear to be far from uncommon in brain-damaged patients: conversely, there is little or nothing in the pathology of hysterical amnesia that has not been observed in the organic syndrome. One case reported in the German literature in 1930 aroused great controversy. A young man developed severe and persistent amnesia following accidental carbon monoxide poisoning. His consciousness was virtually restricted to a second or two and no lasting memory traces could apparently be formed. While the original defect of memory may have been largely, if not

wholly, organic, it was sustained thereafter on a hysterical basis. Conversely, a case has been reported in which the diagnosis, originally hysterical amnesia, had to be altered in light of the discovery that the patient had suffered from progressive brain disease. In such cases, organic and psychogenic factors appear to interact to produce complex and atypical symptoms.

Fugue states. The fugue is a condition in which the individual wanders away from his home or place of work for periods of hours, days, or even weeks. One celebrated case was that of the Rev. Ansell Bourne, described by the U.S. psychologist William James. This clergyman wandered away from home for two months and acquired a new identity. On his return, he was found to have no memory of the period of absence, though it was eventually restored under hypnosis. In not all cases, however, is the basis of the fugue so manifestly psychogenic. Indeed, close observation in some instances may reveal minor alterations in consciousness and behaviour that suggest an organic basis, probably epileptic. According to one view, pathological wandering with subsequent amnesia is due to a constellation of factors, among which are a tendency toward periodic depression, history of a broken home in childhood, and predisposition to states of altered consciousness, even in the absence of organic brain lesion. Psychoanalysts, on the other hand, see in the fugue a symbolic escape from severe emotional conflict.

PARAMNESIA AND CONFABULATION

The term paramnesia was introduced by a German psychiatrist, Emil Kraepelin, in 1886 to denote errors of memory. He distinguished three main varieties; one he called simple memory deceptions, as when one remembers as genuine those events imagined or hallucinated in fantasy or dream. This is not uncommon among confused and amnesic people and also occurs in paranoid states. Kraepelin also wrote of associative memory deceptions, as when a person meeting someone for the first time claims to have seen him on previous occasions. This has been re-named reduplicative paramnesia or simply reduplication. Lastly there was identifying paramnesia, in which a novel situation is experienced as duplicating an earlier situation in every detail; this is now known as *déjà vu* or paramnesia *tout court*. The term confabulation denotes the production of false recollections generally.

Déjà vu. The *déjà vu* experience has aroused considerable interest and is occasionally felt by most people, especially in youth or when they are fatigued. It has also found its way into literature, having been well described by, among other creative writers, Shelley, Dickens, Hawthorne, Tolstoy, and Proust. The curious sense of extreme familiarity may be limited to a single sensory system, such as the sense of hearing, but as a rule it is generalized, affecting all aspects of experience including the subject's own actions. As a rule, it passes off within a few seconds or minutes, though its repercussions may persist for some time. For some epileptics, however, *déjà vu* may continue for hours or even days and can provide a fertile subsoil for delusional elaboration.

In view of its occurrence among organically healthy individuals, *déjà vu* commonly has been regarded as psychogenic and as having its origin in some partly forgotten memory, fantasy, or dream. This explanation has appealed strongly to psychoanalysts; it also gains support from the finding that an experience very similar to *déjà vu* can be induced in normal people by hypnosis. If a picture is presented to a hypnotized person with the instruction to forget it and then is shown with other pictures when he is awake, the subject may report an intense feeling of familiarity that he is at a loss to justify. The *déjà vu* phenomenon also is attributable to minor neurophysiological abnormality; it is frequent in epilepsy. Indeed, *déjà vu* is accepted as a definite sign of epileptic activity originating in the temporal lobe of the brain and may occur as part of the seizure activity or frequently between convulsions. It seems to be more frequent in cases in which the disorder is in the right temporal lobe and has on occasion been evoked by electrical stimulation of the exposed brain during surgery. Some have been tempted

Possible basis in epilepsy

Relief of amnesia through hypnosis

Hypnotic induction of *déjà vu*

to ascribe it to a dysrhythmic electrical discharge in some region of the temporal lobe that is closely associated with memory function.

Reduplicative paramnesia. Reduplication is observed mainly among acutely confused or severely amnesic people; for example, a patient may say that he has been in one or more hospitals that are very similar to his present location and that all bear the same name. The effect also can be induced by showing the person an object such as a picture and by testing him for recognition of the same picture a few minutes later. He is apt to say that he has seen a similar picture but definitely not the one now being shown. This effect appears to depend on loss of a sense of familiarity and on failure to treat a single object seen on a number of occasions as one and the same. It has been reported that reduplication of this kind is typically associated with confabulation, speech disorder (paraphasia), disorientation, and denial of illness.

Confabulation. Spurious memories or fabrications are very common in psychiatric disorders and may take on an expansive and grandiose character. They may also embody obvious elements from fantasy and dream. At a more realistic level, the production of false memories (confabulation) is best studied among sufferers of Korsakoff's syndrome, for whom consciousness and reasoning remain clear. When asked what he did on the previous day, such a person may give a detailed account of a typical day in his life several months or years earlier. Evidently his retrograde amnesia and his disorientation in time provide fertile soil for false reminiscence. When the confabulation embodies dramatic, fanciful elements, it is the exception rather than the rule.

Confabulation once was regarded as one's reaction to the social embarrassment produced by a memory defect—*i.e.*, as an attempt to fill memory gaps plausibly. Despite this possibility, many severely amnesic patients confabulate little, if at all; and there appears to be no relation between the severity of amnesia and frequency of confabulation. In consequence, individual differences in preamnesic personality have been stressed, particularly in regard to suggestibility. While many patients who confabulate are obviously highly suggestible, precise tests of suggestibility have not been used in most clinical evaluations. It also has been claimed that the superficially sociable, but basically secretive, individual is particularly prone to confabulate. The most critical factor appears to be the sufferer's degree of insight into his disorder: it has been observed that the amnesia sufferer who most strongly denies any lapse in memory is most prone to confabulate. By contrast, it also has been claimed that in chronic Korsakoff states the individual's insight into his condition is no guarantee of freedom from confabulation.

While confabulation is pathological by definition, all people include an inventive (and thus spurious) element in their remembering. Indeed, it seems valid to say that all remembering depends heavily on reconstruction rather than on mere reproduction alone. Among amnesiacs, reconstruction is especially drastic, inventive, and error-prone, particularly in regard to chronological sequence. The difference, therefore, between normal and grossly amnesic confabulation may well be one of degree rather than kind.

HYPERMNESIA

Enhancement of memory function (hypermnesia) under hypnosis and in some pathological states was frequently described by 19th-century medical writers; for example, cases were recorded of delirious people who would speak fluently in a language they had not had occasion to use for up to 50 or more years and apparently had forgotten. It was then categorically claimed that anyone under hypnosis would recollect events with invariably greater efficiency than in the waking state. It is true that experience inaccessible to ordinary recall sometimes can be recollected under hypnosis; some have attributed this effect to release from emotional inhibition. Nevertheless, evidence indicates that previously memorized material (*e.g.*, poetry) in many cases is reproduced no better under hypnosis than in the waking state.

Memory prodigies. Few individuals who exhibit excep-

tional memory have been studied extensively. The case of a Russian mnemonist (memory artist), "S," was studied over a period of 30 years, and his story has been delightfully written by a Soviet psychologist (see *Bibliography*). This man's exceptional mnemonic ability seemed largely to depend on an outstandingly vivid, detailed, and persistent visual memory, almost certainly eidetic ("photographic") in nature. "S" also reported an unusual degree of synesthesia, though whether this helped or hindered his feats of memory is not clear. (A person shows signs of synesthesia when he reports that stimulation through one sense leads to experiences in another sense; for example, such a person may say that he sees vivid flashes of colour when he hears music.) Although "S's" highly developed power of concrete visualization made possible feats of memory far beyond the ordinary, he exhibited weakness in abstract thinking.

Exceptional memory capacity is occasionally observed among mathematicians and others with exceptional talent for lightning calculation. A mathematics professor at the University of Edinburgh, for example, was reported to be capable of remarkable feats of long-term memory for personal experiences, music, and verbal material in either English or Latin. This talented mathematician has been said to recall with complete accuracy a list of 25 unrelated words after only a brief effort to memorize, and to recite the value of pi (an endless number) to a thousand places or more. Likewise, some composers and musicians appear to possess exceptional auditory memory, though no systematic study of their attainments appears to have been made. The anatomical or physiological basis of hypermnesia remains most incompletely understood. (O.L.Z.)

BIBLIOGRAPHY

The nature of memory. JACK A. ADAMS, *Human Memory* (1967), provides a well-considered account of memory as viewed from laboratory findings. A somewhat more abbreviated account is JOHN JUNG, *Verbal Learning* (1968). MORRIS MOSCOVITCH (ed.), *Infant Memory: Its Relation to Normal and Pathological Memory in Humans and Other Animals* (1984), is an advanced book for researchers and clinicians. GARY LYNCH, JAMES L. MCGAUGH, and NORMAN M. WEINBERGER (eds.), *Neurobiology of Learning and Memory* (1984), involves mostly animal studies. ALAN BADDELEY, *Human Memory: Theory and Practice* (1990), harmonizes laboratory studies with actual data from brain-damaged patients and is of value to advanced researchers and undergraduates alike. Two collections of conference papers that discuss memory from the connectionist viewpoint are GEOFFREY E. HINTON and JAMES A. ANDERSON (eds.), *Parallel Models of Associative Memory*, updated ed. (1989); and R.G.M. MORRIS (ed.), *Parallel Distributed Processing* (1989). DAVID E. RUMELHART *et al.*, *Parallel Distributed Processing*, 2 vol. (1986), is a landmark collection of articles on the connectionist theory of learning and memory. The origins of connectionist theory are traced in H. CHRISTOPHER LONGUET-HIGGINS, *Mental Processes* (1987), in the section on memory. (B.J.U./Ed.)

Abnormalities of memory. A.R. LURIA, *The Mind of a Mnemonist* (1968, reprinted 1987), is a fascinating account of a "memory prodigy" studied over many years by an outstanding Soviet psychologist. THÉODOLE ARMAND RIBOT, *Diseases of Memory* (1882, reprinted 1977; originally published in French, 1881), is the classical text on disorders of memory. GEORGE A. TALLAND, *Disorders of Memory and Learning* (1968), is a popular survey of memory and some of its disorders. C.W.M. WHITTY and O.L. ZANGWILL (eds.), *Amnesia*, 2nd ed. (1977), considers amnesia from the neurological point of view. DAVID S. OLTON, ELKAN GAMZU, and SUZANNE CORNIN (eds.), *Memory Dysjunctions: An Integration of Animal and Human Research from Preclinical and Clinical Perspectives* (1985), covers a wide variety of topics. ANDREW R. MAYES, *Human Organic Memory Disorders* (1988), is a highly technical and comprehensive book for psychiatric professionals and students. TAKEHIKO YANAGIHARA and RONALD C. PETERSEN (eds.), *Memory Disorders: Research and Clinical Practice* (1991), an advanced text for memory specialists and students, takes a contemporary approach to the clinical assessment of memory and the study of memory dysfunction, making the connection between specific memory functions and specific neuroanatomical and biochemical structures. The memory loss associated with Alzheimer's disease is discussed in DONNA COHEN and CARL EISDORFER, *The Loss of Self* (1986), a practical resource for families and caregivers; and ANTHONY F. JORM, *A Guide to the Understanding of Alzheimer's Disease and Related Disorders* (1987), an overview. (O.L.Z./Ed.)

Photo-graphic memory

Inventive memory

Mental Disorders and Their Treatment

A mental disorder is any illness with significant psychological or behavioral manifestations that is associated with either a painful or distressing symptom or impairment in one or more important areas of functioning. These disorders, in particular their consequences and their treatment, are of more concern and receive more attention now than in the past. Mental disorders have become a more prominent subject of attention for several reasons. They have always been common, but, with the eradication or successful treatment of many of the serious physical illnesses that formerly afflicted humans, mental illness has become a more noticeable cause of suffering and accounts for a higher proportion of those disabled by disease. Moreover, the public has come to expect the medical profession to help it obtain an improved quality of life in mental as well as physical functioning. And, indeed, there has been a proliferation of both pharmacological and psychotherapeutic treatments in this regard. The transfer of many psychiatric patients, some still showing conspicuous symptoms, from mental hospitals into the community has also increased the public's awareness of the importance and prevalence of mental illness.

There is no simple definition of mental disorder that is universally satisfactory. This is partly because mental states or behaviour that are viewed as abnormal or unacceptable in one culture may be regarded as normal or acceptable in another, and in any case it is difficult to draw a line clearly demarcating healthy from abnormal mental functioning.

A narrow definition of mental illness would insist upon the presence of organic disease of the brain, either structural or biochemical. An overly broad definition would define mental illness as simply being the lack or absence of mental health—that is to say, a condition of mental well-being, balance, and resilience in which the individual can successfully work and function and in which he can both withstand and learn to cope with the conflicts and stresses encountered in life. A more generally useful definition than either of the above is that a mental disorder is an illness with significant psychological or behavioral manifestations that occurs in an individual and that is associated either with a painful or distressing symptom, with impairment in one or more important areas of functioning, or with both. The mental disorder may be due to either a psychological, social, biochemical, or genetic dysfunction or disturbance—or a combination of these factors—in the individual.

A mental illness can have an effect on every aspect of a

person's life, including thinking, feeling, mood, and outlook and such areas of external activity as family and marital life, sexual activity, work, recreation, and management of material affairs. Most mental disorders negatively affect how individuals feel about themselves and impair their capacity for participating in mutually rewarding relationships.

Psychopathology is the systematic study of the significant causes, processes, and symptomatic manifestations of mental disorders. The meticulous study, observation, and inquiry that characterize the discipline of psychopathology are, in turn, the basis for the practice of psychiatry (*i.e.*, the science and practice of diagnosing and treating mental disorders, as well as dealing with their prevention), clinical psychology, and counseling. Psychiatry and related disciplines embrace a wide spectrum of techniques and approaches for treating mental illnesses. These include the use of psychoactive drugs to correct biochemical imbalances in the brain or otherwise to relieve depression, anxiety, and other painful emotional states.

Another important group of treatments comprises the psychotherapies, which seek to treat mental disorders by psychological means and which involve verbal communication between the patient and a trained person in the context of a therapeutic interpersonal relationship between them. An important variant of this latter mode of treatment is behaviour therapy, which concentrates on changing or modifying observable pathological behaviours by the use of conditioning and other experimentally derived principles of learning.

This article discusses the types, causes, and treatment of mental disorders. Neurological diseases with behavioral manifestations are treated in *NERVES AND NERVOUS SYSTEMS*. Alcoholism and other substance use disorders are discussed in *ALCOHOL AND DRUG CONSUMPTION*. Disorders of sexual functioning and behaviour are treated in *SEX AND SEXUALITY*. Mental retardation is treated in *INTELLIGENCE, HUMAN*. Tests used to evaluate mental health and functioning are discussed in *PSYCHOLOGICAL TESTS AND MEASUREMENT*. The various theories of personality structure and dynamics are treated in *PERSONALITY*, while human emotion and motivation are discussed in *EMOTION, HUMAN*.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia* section 436, and the *Index*.

This article is divided into the following sections:

Types and causes of mental disorders 842

Classification and epidemiology 842

Classification

Epidemiology

Theories of causation 843

Organic and hereditary etiologies

Psychodynamic etiologies

Behavioral etiology

Other etiologies

Major diagnostic categories 845

Organic mental disorders

Schizophrenia

Paranoid disorders

Mood disorders

Anxiety disorders

Obsessive-compulsive disorder

Posttraumatic stress disorder

Somatiform disorders

Dissociative disorders

Personality disorders

Psychosexual disorders

Disorders usually first evident in infancy, childhood, or adolescence

Other mental disorders

Treatment of mental disorders 851

Historical overview 851

Early history

The mental hospital era

The biological movement

Development of psychotherapy

Development of physical and pharmacological treatments

Deinstitutionalization

Development of behaviour therapy

The mental health profession in the late 20th century

Physiological treatments 853

Pharmacological treatments

Electroconvulsive treatment

Psychosurgery

The psychotherapies 856

Dynamic psychotherapies

Behavioral psychotherapy

Other therapies

Bibliography 859

Types and causes of mental disorders

CLASSIFICATION AND EPIDEMIOLOGY

Psychiatric classification attempts to bring order to the enormous diversity of mental symptoms, syndromes, and illnesses that are encountered in clinical practice. Epidemiology is the measurement of the prevalence, or frequency of occurrence, of these psychiatric disorders in different human populations.

Classification. Diagnosis is the process of identifying an illness by studying its signs and symptoms and by considering the patient's history. Much of this information is gathered by the mental health practitioner (e.g., psychiatrist, psychologist, social worker, or counselor) during initial interviews with the patient, who describes the main complaints and symptoms and any past ones and briefly gives a personal history and current situation. The practitioner may administer any of several psychological tests to the patient and may supplement these with a physical and a neurological examination. These data, along with the practitioner's own observations of the patient and of the patient's interaction with him, form the basis for a preliminary diagnostic assessment. For the practitioner, diagnosis involves finding the most prominent or significant symptoms, upon which the patient's disorder can be assigned to a category as a first stage toward treatment.

Classification systems in psychiatry aim to distinguish groups of patients who share the same or related clinical symptoms in order to provide an appropriate therapy and accurately predict the prospects of recovery for any individual member of that group. Thus, the diagnosis of, for example, depression having been made, it becomes logical to consider antidepressant drugs when preparing a course of treatment.

The diagnostic terms of psychiatry have been introduced at various stages of the discipline's development and from very different theoretical standpoints. Sometimes two words with quite different derivations have come to mean almost the same thing—for example, *dementia praecox* and *schizophrenia*. Sometimes a word, such as *hysteria*, carries many different meanings, depending on the psychiatrist's theoretical orientation.

Psychiatry is hampered by the fact that the cause of many mental illnesses is unknown, and so convenient diagnostic distinctions cannot be made between such illnesses, as they can, for instance, in infectious medicine, where infection with a specific type of bacterium is a reliable indicator for a diagnosis of tuberculosis. But the greatest difficulty presented by mental disorders as far as classification and diagnosis are concerned is that the same symptoms are often found in patients with different or unrelated disorders, or a patient may show a mix of symptoms properly belonging to several different disorders. Thus, although the categories of mental illness are defined according to symptom patterns, course, and outcome, the illnesses of many patients constitute intermediate cases between such categories, and the categories themselves may not necessarily represent distinct disease entities and are often poorly defined.

The two most frequently used systems of psychiatric classification are the *International Classification of Diseases*, produced by the World Health Organization, and the *Diagnostic and Statistical Manual of Mental Disorders* (DSM), produced by the American Psychiatric Association. The 10th revision of the former, published in 1992, is widely used in western Europe and other parts of the world for epidemiological and administrative purposes. Its nomenclature is deliberately conservative in conception so that it can be used by clinicians and mental health care systems in different countries.

This article, however, will follow the fourth edition of the DSM, which was published in 1994. The DSM differs from the International Classification in its introduction of precisely described criteria for each diagnostic category; its categorizations are usually based upon the detailed description of symptoms.

The DSM-IV has been widely used, especially in the United States, and its detailed descriptions of diagnostic criteria have been useful in eradicating the inconsistencies of earlier classifications. However, there are still some major

problems in its everyday clinical use. Chief among these problems is the DSM's innovative and controversial abandonment of the general categories of psychosis and neurosis in its classificatory scheme. These terms have been and still are widely used to distinguish between classes of mental disorders, though there are various mental illnesses, such as personality disorders, that cannot be classified as either psychoses or neuroses.

Psychoses. Psychoses are major mental illnesses that are characterized by severe symptoms such as delusions, hallucinations, disturbances of the thinking process, and defects of judgment and insight. Persons with psychoses exhibit a disturbance or disorganization of thought, emotion, and behaviour so profound that they are often unable to function in everyday life and may be incapacitated or disabled. Such individuals are often unable to realize that their subjective perceptions and feelings do not correlate with objective reality, a phenomenon evinced by persons with psychoses who do not know or will not believe that they are ill despite the distress they feel and their obvious confusion concerning the outside world. Traditionally, the psychoses have been broadly divided into organic and functional psychoses. Organic psychoses were believed to result from a physical defect of or damage to the brain. Functional psychoses were believed to have no physical brain disease evident upon clinical examination. Much recent research suggests that this distinction between organic and functional is probably inaccurate. Most psychoses are now believed to result from some structural or biochemical change in the brain.

Neuroses. Neuroses, or psychoneuroses, are less serious disorders in which a person may experience such negative feelings as anxiety or depression. His functioning may be significantly impaired, but his personality remains relatively intact, he maintains a capacity for recognizing and objectively evaluating reality, and he is basically able to function in everyday life. In contrast to a person with a psychosis, an individual with a neurosis knows or can be made to realize that he is ill, and he usually wants to get well and return to a normal state. His chances for recovery are better than those of someone with a psychosis. The symptoms of neurosis may sometimes resemble the coping mechanisms used in everyday life by most people, but in persons with a neurosis these defensive reactions are inappropriately severe or prolonged in response to an external stress. Anxiety disorders, phobic disorder, conversion disorder, obsessive-compulsive disorder, and depressive disorders have been traditionally classified as neuroses.

Epidemiology. Epidemiology is the study of the distribution of disease in different populations. Prevalence denotes the number of cases of a condition present at a particular time or over a specified period, while incidence denotes the number of new cases occurring in a defined time period. Epidemiology is also concerned with the social, economic, or other contexts in which mental illnesses arise.

The understanding of mental disorders is aided by knowledge of the rate and frequency with which they occur in different societies and cultures. Looking at the worldwide prevalence of mental disorders reveals many surprising findings. It is remarkable, for instance, how constant the rate of developing schizophrenia is: in widely different cultures there is a lifetime risk of developing the illness of approximately just under 1 percent.

Gradual historical changes in the incidence and prevalence of particular disorders have often been described, but it is very difficult to obtain firm evidence that such changes have actually occurred. On the other hand, prevalence has been seen to increase for a few syndromes owing to general changes in living conditions over time. For example, dementia inevitably develops in some 20 percent of those persons over age 80, so that with the increase in life expectancy common to developed countries the number of people with dementia is bound to increase. Other factors, such as the presence of small quantities of aluminum in drinking water, may also play a part in the increased prevalence of dementia. There also seems to be some evidence of an increased prevalence of mood disorders over the past century.

Several large-scale epidemiological studies have been

Mixed symptoms in many mental disorders

Responsiveness of neuroses to therapy

Prevalence of mental disorders

conducted to determine the incidence and prevalence of mental disorders in the general population. Simple statistics based on those people actually under treatment for mental disorders cannot be relied upon in making such a determination because the number of those who have sought treatment is substantially smaller than the actual number of people afflicted with mental disorders, many of whom do not seek professional treatment. Moreover, surveys to determine incidence and prevalence depend for their statistics on the clinical judgment of the survey takers, which can always be fallible because there are no objective tests for the assessment of mental illness. Given such objections, one ambitious study conducted by the National Institutes of Mental Health in the United States examined thousands of persons in several American localities and yielded the following results concerning the prevalence of mental disorders in the general population. About 1.1 percent of those surveyed were found to have schizophrenia; 9.5 percent had depression; and 13.3 percent had phobias or other anxiety disorders.

There is a relatively strong epidemiological association between socioeconomic class and the occurrence of certain types of mental disorders and of general patterns of mental health. One study found that the lower the socioeconomic class, the greater the prevalence of psychotic disorders; schizophrenia was found to be 11 times more frequent among the lowest of the five classes surveyed (unskilled manual workers) than among the highest class (professionals). (Anxiety disorders were found to be more common among the middle class, however.) Two possible explanations for the elevated frequency of schizophrenia among the poor would be that persons with schizophrenia "drift downward" to the lowest socioeconomic class because they are impaired by their illness, or alternatively that unfavourable sociocultural conditions create circumstances that help induce the illness.

The manifestation of particular psychiatric symptoms is sometimes closely associated with particular epochs or periods in life. The symptoms of autism are usually evident by early childhood, for example. Childhood and adolescence may produce a variety of psychiatric symptoms peculiar to those periods of life. Anorexia nervosa, several types of schizophrenia, drug abuse, and bipolar disorder often first appear during adolescence or in young adult life. Alcohol dependence and its consequences, paranoid schizophrenia, and repeated attacks of depression are more likely to occur in middle age. Involitional melancholia and presenile dementias typically occur in late middle age, while senile and arteriosclerotic dementias are characteristic of the elderly.

There are also marked sex differences in the incidence of certain types of mental illness. For instance, anorexia nervosa is 20 times more common in girls than boys; men tend to develop schizophrenia at a younger age than women; depression is more common in women than men; and many sexual deviations occur almost exclusively in men.

THEORIES OF CAUSATION

Very often the etiology, or cause, of a particular type of mental disorder is unknown or is understood only to a very limited extent. The situation is complicated by the fact that a mental disorder such as schizophrenia may be caused by a combination and interaction of several factors, including a probable genetic predisposition to developing the disease, a postulated biochemical imbalance in the brain, and a cluster of stressful life events that help to precipitate the actual onset of the illness. The predominance of these and other factors probably varies from person to person in schizophrenia. A similarly complex interaction of constitutional, developmental, and social factors can influence the formation of mood and anxiety disorders.

No single theory of causation can explain all mental disorders or even all those of a particular type, and, moreover, the same type of disorder may have different causes in different persons; e.g., an obsessive-compulsive disorder may have its origins in a biochemical imbalance, in an unconscious emotional conflict, in faulty learning processes, or in a combination of these. The fact that quite different therapeutic approaches can produce equal improvements in

different patients with the same type of disorder underscores the complex and ambiguous nature of the causes of mental illness. The major theoretical and research approaches to the causation of mental disorders are treated below.

Organic and hereditary etiologies. Organic explanations of mental illness have usually been genetic, biochemical, neuropathological, or a combination of these.

Genetics. The study of the genetic causes of mental disorders involves both the laboratory analysis of the human genome and the statistical analysis of the frequency of a particular disorder's occurrence among individuals who share related genes—i.e., family members and particularly twins. Family-risk studies compare the observed frequency of occurrence of a mental illness in close relatives of the patient with its frequency in the general population. First-degree relatives (parents, siblings, and children) share 50 percent of their genetic material with the patient, and higher rates of the illness in these relatives than expected indicate a possible genetic factor. In twin studies the frequency of occurrence of the illness in both members of pairs of identical (monozygotic) twins is compared with its frequency in both members of a pair of fraternal (dizygotic) twins. A higher concordance for disease among the identical than the fraternal twins suggests a genetic component. Further information on the relative importance of genetic and environmental factors accrues from comparing identical twins reared together with those reared apart. Adoption studies comparing adopted children whose biological parents had the illness with those whose parents did not can also be useful in separating biological from environmental influences.

Such studies have demonstrated a clear role for genetic factors in the causation of schizophrenia. When one parent is found to have the disorder, the probability that his children will develop schizophrenia is at least 10 times higher (about a 12 percent risk probability) than it is for children in the general population (about a 1 percent risk probability). If both parents have schizophrenia, their children stand anywhere from a 35 to a 65 percent probability of developing schizophrenia. If one member of a pair of fraternal twins develops schizophrenia, there is about a 12 percent chance that the other twin will also develop the disorder. If one member of a pair of identical twins has schizophrenia, the other identical twin has at least a 40–50 percent chance of developing the disease. Genetic factors seem to play a less significant role in the causation of other psychotic disorders and in personality disorders. However, studies have demonstrated a probable role for genetic factors in the causation of many mood disorders and some anxiety disorders.

Biochemistry. If a mental disease is caused by a biochemical abnormality, investigation of the brain at the site where the biochemical imbalance occurs should show neurochemical differences from normal. In practice such a simplistic approach is fraught with practical, methodological, and ethical difficulties. The living human brain is not readily accessible to direct investigation, and the dead brain undergoes chemical change; moreover, findings of abnormalities in cerebrospinal fluid, blood, or urine may have no relevance to the question of a presumed biochemical imbalance in the brain. It is difficult to study human mental illnesses by using animals as analogs, because most mental disorders either do not occur or are not recognizable in animals. Even when biochemical abnormalities have been found in persons with mental disorders, it is difficult to know whether such abnormalities are the cause or the result of the illness, of its treatment, or of other consequences. Despite these problems, progress has been made in unraveling the biochemistry of mood disorders, schizophrenia, and some of the dementias.

Certain drugs have been demonstrated to have beneficial effects upon mental illnesses. Antidepressant, antipsychotic, and anti-anxiety drugs are thought to achieve their therapeutic results by the selective inhibition or enhancement of the quantities, action, or breakdown of neurotransmitters in the brain. Neurotransmitters are a group of chemical agents that are released by neurons (nerve cells) to stimulate neighbouring neurons, which thus allows im-

Genetic predisposition to mental disorders

Role of neurotransmitters

pulses to be passed from one cell to the next throughout the nervous system. Neurotransmitters play a key role in transmitting nerve impulses across the microscopic gap (synaptic cleft) that exists between neurons. The release of such neurotransmitters is stimulated by the electrical activity of the cell. Norepinephrine, dopamine, acetylcholine, and serotonin are among the principal neurotransmitters. Some neurotransmitters excite or activate neurons, while others act as inhibiting substances. Abnormally low or high concentrations of neurotransmitters at sites in the brain are thought to change the synaptic activities of neurons and thus ultimately lead to the disturbances of mood, emotion, or thought found in various mental disorders.

Neuropathology. In the past the pathological study of the brain at postmortem revealed information upon which great advances in understanding of the etiology of neurological and some mental disorders were based, which led to the German psychiatrist Wilhelm Griesinger's postulate: "All mental illness is disease of the brain." The application of the principles of pathology to general paresis of the insane, one of the most common conditions found in mental hospitals in the late 19th century, resulted in the discovery that this was a form of neurosyphilis and was caused by infection with the spirochete bacterium *Treponema pallidum*. The examination of the brains of patients with other forms of dementia has given useful information concerning other causes of this syndrome—for example, Alzheimer's disease and arteriosclerosis. The pinpointing of abnormalities of specific areas of the brain has aided understanding of some abnormal mental functions, such as disturbances of memory or speech disorders. Recent advances in neuroimaging techniques have expanded the ability to investigate brain abnormalities in patients with a wide variety of mental illnesses and no longer require that these studies occur postmortem.

Psychodynamic etiologies. Up to the 1950s theories of the etiology of mental disorders, especially of neuroses and personality disorders, were dominated in the United States by Freudian psychoanalysis and the derivative theories of the post-Freudians. In western Europe the influence of Freudian theory upon psychiatric theory diminished after World War II.

Theories of personality development. Freudian and other psychodynamic theories view neurotic symptoms as arising from intrapsychic conflict—that is, as being caused by conflicting motives, drives, impulses, and feelings held within various components of the mind. Central to psychoanalytic theory is the postulated existence of the unconscious, which is that part of the mind whose processes and functions are inaccessible to the individual's conscious awareness or scrutiny. One of the functions of the unconscious is thought to be that of a repository for traumatic memories, feelings, ideas, wishes, and drives that are threatening, abhorrent, anxiety-provoking, or socially or ethically unacceptable to the individual. These mental contents may at some time be pushed out of conscious awareness but remain actively held in the unconscious. This process is a defense mechanism for protecting the individual from the anxiety or other psychic pain associated with those contents and is known as repression. The repressed mental contents held in the unconscious retain much of the psychic energy or power that was originally attached to them, however, and they can continue to influence significantly the mental life of the individual even though (or because) he is no longer aware of them.

The natural tendency for repressed drives or feelings, according to this theory, is to break through into conscious awareness so that the individual can seek the gratification, fulfillment, or resolution of them. But this threatened release of forbidden impulses or memories provokes anxiety and is seen as threatening, and a variety of defense mechanisms may then come into play to provide relief from the state of psychic conflict. Through reaction formation, projection, regression, sublimation, rationalization, and other defense mechanisms, some component of the unwelcome mental contents can emerge into consciousness in a disguised or attenuated form and thus provide partial relief to the individual. Later, perhaps in adult life, some event or situation in the person's life triggers the abnormal dis-

charge of the dammed-up or strangulated emotional energy in the form of neurotic symptoms in a manner mediated by defense mechanisms. Such symptoms can form the basis of neurotic disorders such as conversion and somatoform disorders, anxiety disorders, obsessional disorders, and depressive disorders. Since the symptoms represent a compromise within the mind between letting the repressed mental contents out and continuing to deny all conscious knowledge of them, the particular character and aspects of an individual's symptoms and neurotic concerns bear an inner meaning that symbolically represents the underlying intrapsychic conflict. Psychoanalysis and other dynamic therapies associate the patient's controlled and therapeutic recovery to conscious awareness of repressed mental conflicts, and his understanding of their influence on both his past history and his present difficulties, with the relief of symptoms and improved mental functioning.

Freudian theory views childhood as the primary breeding ground of neurotic conflicts. This is because children are relatively helpless and are dependent on their parents for love, care, security, and support and because their psychosexual, aggressive, and other impulses are not yet integrated into a stable personality framework. Children are thus subject to emotional traumas, deprivations, and frustrations which they lack the resources to cope with and which can become grounds for intrapsychic conflicts that are not resolved but rather merely held in abeyance through repression, producing insecurity, unease, or guilt and subtly influencing the individual's developing personality, interests, attitudes, and ability to cope with later stresses.

Non-Freudian psychodynamics. Psychoanalytic theory's emphasis on the unconscious mind and its influence on human behaviour resulted in a proliferation of other, related theories of causation incorporating many basic psychoanalytic precepts. Most subsequent psychotherapies have stressed in their theories of causation aspects of earlier, maladaptive psychological development that had been missed or underemphasized by orthodox psychoanalysis, or they have incorporated insights taken from learning theory. The Swiss psychiatrist Carl Jung, for instance, concentrated on the individual's need for spiritual development and concluded that neurotic symptoms could arise from a lack of self-fulfillment in this regard. The Austrian psychiatrist Alfred Adler emphasized the importance of feelings of inferiority and the unsatisfactory attempts to compensate for it as important causes of neurosis. Neo-Freudian authorities such as Harry Stack Sullivan, Karen Horney, and Erich Fromm modified Freudian theory by emphasizing social relationships and cultural and environmental factors as being important in the formation of mental disorders. Many other highly specific theories of causation have been developed by particular psychotherapies, and, in general, psychiatric scrutiny has come to extend far beyond the confines of early psychosexual development that were originally posited by Freud as the prime ground for the causation of neuroses. More modern psychodynamic theories have moved away from the idea of explaining and treating neurosis on the basis of a defect in a single psychological system and have instead adopted a more complex notion of multiple causes, including emotional, psychosexual, social, cultural, and existential ones. A notable trend in the more recently developed psychotherapies has been the incorporation of approaches derived from theories of learning. Such psychotherapies pay special attention to acquired faulty mental processes and maladaptive behavioral responses that act to sustain neurotic symptoms, and there has generally been increased interest in the patient's present circumstances and his learned responses to those conditions as a causative factor in mental illness. In this way, psychoanalytic theory and behavioral theory have tended somewhat to converge and intermingle in their views of disease causation.

Behavioral etiology. Behavioral theories for the causation of mental disorders, especially neurotic symptoms, are based upon learning theory, which was in turn largely derived from the study of the behaviour of animals in laboratory settings. Most important theories in this area arose out of the work of the Russian physiologist Ivan Pavlov and such American psychologists as Edward L. Thorndike,

Childhood
experience
and
neurosis

Clark L. Hull, John B. Watson, Edward C. Tolman, and B.F. Skinner. In the classical Pavlovian model of conditioning, an unconditioned stimulus is followed by an appropriate response; for example, food placed in a dog's mouth is followed by the dog's salivating. If a bell is rung just before food is offered to a dog, eventually the dog will salivate at the sound of the bell only, even though no food is offered. Because the bell could not originally evoke salivation in the dog (and hence was a neutral stimulus) but came to evoke salivation because it was repeatedly paired with the offering of food, it is called a conditioned stimulus. The dog's salivation at the sound of the bell alone is called a conditioned response. If the conditioned stimulus (the bell) is no longer paired with the unconditioned stimulus (the food), extinction of the conditioned response gradually occurs (the dog ceases to salivate at the sound of the bell alone).

Behavioral theories for the causation of mental disorders rest largely upon the assumption that the symptoms or symptomatic behaviour found in persons with various neuroses (particularly phobias and other anxiety disorders) can be regarded as learned behaviours that have been built up into conditioned responses. In the case of phobias, for example, a person who has once been exposed to an inherently frightening situation afterward experiences anxiety even at neutral objects that were merely associated with that situation at the time but that should not reasonably produce anxiety; *e.g.*, a child who has had a frightful experience with a bird may subsequently have an unreasonable dread of feathers. The neutral object alone is enough to arouse anxiety, and the person's subsequent effort to avoid that object is a learned behavioral response that is self-reinforcing, since the person does indeed procure a reduction of his anxiety by avoiding the feared object and is thus likely to continue to avoid it in the future. It is only by confronting the object that he can eventually lose his irrational association-based fear of it.

Other etiologies. Mental illness may be contagious in a psychological sense; that is, close contact with an individual who has symptoms may result in the transmission of those symptoms to one or many others who were previously unaffected. This may occur either through the powerful influence of long-term cohabitation of one person with one other—a phenomenon known as *folie à deux*—or through the volatile collective emotions of a group—mass hysteria.

Social values can sometimes determine or encourage the formation of particular syndromes. Prime examples of this are anorexia nervosa and bulimia nervosa, which predominantly affect young white females in affluent Western societies. The value and attractiveness of physical thinness are communicated via the media, and an eating disorder resulting in emaciation subsequently occurs in some susceptible individuals.

Another approach to the causation of mental disorders focuses on the effects and consequences of stress, which is a state of bodily or mental tension resulting from external factors such as marital conflicts, excessive work demands, or serious financial problems. Stress may cause psychosomatic illnesses or an exacerbation or worsening of real somatic illnesses. An accumulation of stressful life events can contribute to the development of depression in psychologically vulnerable individuals.

MAJOR DIAGNOSTIC CATEGORIES

Organic mental disorders. This category includes both those psychological or behavioral abnormalities that arise from structural disease of the brain and those that arise from brain dysfunction caused by disease outside the brain. These conditions differ from those of other mental illnesses in that they have a definite and ascertainable cause—*i.e.*, brain disease. However, the importance of the distinction (between organic and functional) has become less clear as research has demonstrated that brain abnormalities are associated with many psychiatric illnesses. Treatment, when possible, is aimed at both the symptoms and the underlying physical dysfunction in the brain.

There are several types of psychiatric syndromes that arise clearly from organic brain disease, chief among them being

dementia and delirium. Dementia is a gradual and progressive loss of such intellectual abilities as thinking, remembering, paying attention, judging, and perceiving, without an accompanying disturbance of consciousness. The syndrome may also be marked by the onset of personality changes. Dementia is usually a chronic condition and frequently worsens over the long term. Delirium is a diffuse or generalized intellectual impairment marked by a clouded or confused state of consciousness, an inability to attend to one's surroundings, difficulty in thinking coherently, a tendency to perceptual disturbances such as hallucinations, and difficulty in sleeping. Delirium is generally an acute condition. Amnesia (*i.e.*, a gross loss of recent memory and time sense without other intellectual impairment) is another specific psychological impairment associated with organic brain disease.

In the diagnosis of suspected organic disorders, a full history must first be taken of the patient and his mental state must be examined in detail, with additional tests for particular functions added if necessary. A physical examination is also performed with special attention to the central nervous system. In order to determine whether a metabolic or other biochemical imbalance is causing the condition, blood and urine tests, liver function tests, thyroid function tests, and other evaluations may be performed. Chest and skull X rays may be taken, and computed tomography (CT) scanning or magnetic resonance imaging (MRI) is frequently used to reveal focal or generalized brain disease. Electroencephalography (EEG) may show localized abnormalities in the electrical conduction of the brain caused by a lesion. Detailed psychological testing may reveal more specific perceptual, memory, or other disabilities.

Senile and presenile dementia. In these dementias there is a progressive intellectual impairment that proceeds to lethargy, inactivity, and gross physical deterioration and eventually to death within a few years. Presenile dementias are arbitrarily defined as those that begin in persons under the age of 65. In old age the most common causes of dementia are Alzheimer's disease and cerebral arteriosclerosis. Dementia from Alzheimer's disease usually begins in people over age 65 and is more common in women than in men. It begins with incidences of forgetfulness, which become more frequent and serious, and the disturbances of memory, personality, and mood progress steadily toward physical deterioration and death within a few years. In dementia caused by cerebral arteriosclerosis, multiple areas of destruction of the brain (infarcts) are caused by pieces of the damaged arteries outside the skull lodging in the small arteries of the brain. The course of the illness is rapid, with periods of deterioration followed by periods of slight improvement. Death may be delayed slightly longer than with dementia from Alzheimer's disease and often occurs from ischemic heart disease, causing a heart attack, or from massive cerebral infarction, causing a stroke.

Other causes of dementia include Pick's disease, a rare inherited condition that occurs in women twice as often as men, usually between the ages of 50 and 60; Huntington's chorea, an inherited disease that usually begins at about the age of 40 with involuntary movements and proceeds to dementia and death within 15 years; and Creutzfeldt-Jakob disease, a rare condition that is caused by an abnormal protein found in the brain called a prion. Head injury—for instance, resulting from a boxing career or from an accident—may produce dementia. Infection—for example, with syphilis or encephalitis—various tumours, toxic conditions such as chronic alcoholism or heavy metal poisoning, metabolic illnesses such as liver failure, reduced oxygen to the brain due to anemia or carbon monoxide poisoning, and the inadequate intake or metabolism of certain vitamins may all result in dementia.

There is no specific treatment for the symptoms of dementia; the underlying physical cause needs to be identified and treated when possible. The goals of care of the individual with dementia are to relieve distress, prevent behaviour that might result in injury, and optimize remaining physical and psychological faculties.

Substance-induced disorders. A variety of psychiatric conditions can result from the use of alcohol or other drugs. Mental disorders resulting from the ingestion of al-

The role of conditioned responses

Alzheimer's disease

cohol include intoxication, withdrawal, hallucinations, and amnesia. Similar syndromes may occur following the use of other drugs that affect the central nervous system (see ALCOHOL AND DRUG CONSUMPTION). Those drugs most commonly used recreationally to alter mood are barbiturates, opioids (e.g., heroin), cocaine, amphetamines, hallucinogens such as LSD (lysergic acid diethylamide), marijuana, and tobacco. Treatment is directed at alleviating symptoms and preventing further abuse of the substance. Substance abuse implies a sustained pattern of use, resulting in an impairment of the individual's daily functioning. Substance dependence occurs when the individual who is abusing a substance becomes physically tolerant of it; in other words, he must administer markedly increased amounts of the drug to achieve the same effect that would have previously occurred with smaller doses. Dependence is also characterized by withdrawal symptoms, which may include tremors, nausea, and anxiety, that follow the cessation of drug use or decreases in dosage.

Schizophrenia. The term *schizophrenia* was introduced by Swiss psychiatrist Eugen Bleuler in 1911 to describe what he considered to be a group of severe mental illnesses with related characteristics; *schizophrenia* eventually replaced the earlier term *dementia praecox*, which the German psychiatrist Emil Kraepelin had first used in 1899 to distinguish the disease from what is now called bipolar disorder. Individuals with schizophrenia exhibit a wide variety of symptoms; thus, although different experts may agree as to whether a particular individual suffers from the condition, they might disagree about which symptoms are essential in clinically defining schizophrenia.

The annual prevalence of schizophrenia—that is, the number of cases both old and new recorded in one year—is between two and four per 1,000 persons. The lifetime risk of developing the illness is between seven and nine per 1,000. Schizophrenia is the single largest cause of admissions to mental hospitals, and it accounts for an even larger proportion of the permanent populations of such institutions. It is a severe and frequently chronic illness that typically first manifests itself during the teen years or during early adult life. More severe levels of impairment and personality disorganization occur in schizophrenia than in almost any other mental disorder.

Clinical features. The principal clinical signs of schizophrenia are delusions, hallucinations, a loosening or incoherence of a person's thought processes and train of associations, deficiencies in feeling appropriate or normal emotions, and a withdrawal from reality. A delusion is a false or irrational belief that is firmly held despite obvious or objective evidence to the contrary. The delusions of individuals with schizophrenia may be persecutory, grandiose, religious, sexual, or hypochondriacal in nature, or they may be concerned with other topics. Delusions of reference, in which the person attributes a special, irrational, and usually negative significance to other people, objects, or events, are common in the disease. Especially characteristic of schizophrenia are delusions in which the individual believes his thinking processes, body parts, or actions or impulses are controlled or dictated by some external force.

Hallucinations are false sensory perceptions that are experienced without an external stimulus but that nevertheless seem real to the person who is experiencing them. Auditory hallucinations, experienced as "voices" and characteristically heard commenting negatively about the affected person in the third person, are prominent in schizophrenia. Hallucinations of touch, taste, smell, and bodily sensation may also occur. Disorders of thinking vary in nature but are quite common in schizophrenia. Thought disorders may consist of a loosening of associations, so that the speaker jumps from one idea or topic to another unrelated one in an illogical, inappropriate, or disorganized way. At its most serious, this incoherence of thought extends into pronunciation itself, and the speaker's words become garbled or unrecognizable. Speech may also be overly concrete and inexpressive; it may be repetitive, or, though voluble, it may convey little or no real information. Usually an individual with schizophrenia has little or no insight into his own condition and realizes neither that he

is suffering from a mental illness nor that his thinking is disordered.

Among the so-called negative symptoms of schizophrenia are a blunting or flattening of the person's ability to experience (or at least to express) emotion, indicated by speaking in a monotone and by a peculiar lack of facial expressions. The person's sense of self (*i.e.*, of who he is) may be disturbed. A person with schizophrenia may be apathetic and may lack the drive and ability to pursue a course of action to its logical conclusion, may withdraw from society, may become detached from others, or may become preoccupied with bizarre or nonsensical fantasies. Such symptoms are more typical of chronic rather than of acute schizophrenia.

Experts have recognized many different types of schizophrenia, and there are intermediate stages between the disease and other conditions. Five major types of schizophrenia are recognized by the DSM-IV: the disorganized type, the catatonic type, the paranoid type, the undifferentiated type, and the residual type. Disorganized schizophrenia is characterized by inappropriate emotional responses and by incoherent thought and speech. Catatonic schizophrenia is marked by striking motor behaviour, such as remaining motionless in a rigid posture for hours or even days, and by stupor or mutism. Paranoid schizophrenia is characterized by the presence of prominent delusions of a persecutory or grandiose nature. The undifferentiated and residual types are marked by the absence of these distinct features; the residual type is, moreover, a less severe diagnosis.

Course and prognosis. The course of schizophrenia is variable. Some individuals with schizophrenia continue to function fairly well and are able to live independently; some have recurrent episodes of the illness with some negative effect on their overall level of function; and some deteriorate into chronic schizophrenia with severe disability. The prognosis for individuals with schizophrenia has improved owing to the use of antipsychotic drugs and community supportive measures.

About 10 percent of individuals with schizophrenia commit suicide. The prognosis of schizophrenia is poorer when it has a gradual rather than a sudden onset, when the affected individual is quite young at the onset of the disease, when there is a long duration of illness, when the individual exhibits blunted feelings or has displayed an abnormal personality previous to the onset of the disease, and when such social factors as never having been married, poor sexual adjustment, a poor employment record, or social isolation exist in the person's history.

Etiology. An enormous amount of research has been performed to try to determine the causes of schizophrenia. Family, twin, and adoption studies provide strong evidence to support an important genetic contribution. A recent study found that children born to men who are over the age of 50 are nearly three times more likely to have schizophrenia than those born to younger men. Stressful life events are known to trigger or quicken the onset of schizophrenia or to cause relapse. Some abnormal neurological signs have been found in individuals with schizophrenia, and it is possible that brain damage, perhaps occurring at birth, may be a cause in some cases. Other studies hypothesize that the cause of schizophrenia is linked to a virus or to the abnormal activity of genes governing the formation of nerve fibres in the brain. Various biochemical abnormalities have been reported in persons with schizophrenia. There is evidence that the abnormal coordination of neurotransmitters such as dopamine, glutamate, and serotonin may be involved in the development of the disease. The neurodevelopmental model currently best incorporates how biological and environmental influences affect the onset and progression of schizophrenia. Functional neuroimaging studies have demonstrated that the frontal and limbic areas of the brain are abnormal in persons with schizophrenia.

Research also has been performed to determine whether the prenatal care used in the families of individuals with schizophrenia contributes to the development of the disease. There has also been extensive interest in such factors as social class, place of residence, migration, and social

Major types of schizophrenia

Hallucinations in schizophrenia

isolation. Neither family dynamics nor social disadvantage has been proved to be a causative agent.

Treatment. The most successful treatment approaches combine the use of medications and supportive therapy. New "atypical" antipsychotic medications such as clozapine, risperidone, and olanzapine have proved effective in relieving or eliminating such symptoms as delusions, hallucinations, thought disorders, agitation, and violent behaviour. These medications also have fewer side effects than the more traditional antipsychotic medications. Long-term maintenance on such medications also reduces the rate of relapse. Psychotherapy may help relieve feelings of helplessness and isolation, reinforce healthy or positive tendencies, allow the affected individual to distinguish between his psychotic perceptions and reality, and explore any underlying emotional conflicts that might be exacerbating the condition. Occupational therapy and regular visits from a social worker or psychiatric nurse may be beneficial. It is sometimes useful to counsel the relatives of individuals with schizophrenia living at home. Support groups for persons with schizophrenia and their family members have become extremely important resources for dealing with the disorder.

Paranoid disorders. Paranoia is a syndrome in which a person thinks or believes, without justification, that other people are plotting or conspiring against, persecuting, harming, or harassing him in some way. Paranoid thinking frequently causes a person to interpret or exaggerate innocuous or trivial incidents in a self-referent way—*e.g.*, to see two people talking at a distance and to assume irrationally that they are plotting against or criticizing him. Grandiosity or delusions of grandeur, which consist of exaggerated and unjustified ideas of a person's own importance, wealth, or power, frequently coexist with the classic persecutory orientation in paranoia. Paranoia or paranoid thinking can be a prominent or primary feature in schizophrenia (paranoid schizophrenia), personality disorders, senile dementias, mood disorders, and bipolar disorder. Persons with paranoid disorder are usually otherwise normal people who may be simply abnormally suspicious, or they may have an unshakable and highly elaborate delusional system involving worldwide conspiracies against them. A special type of paranoia is delusional jealousy, in which a person delusionally believes or suspects that his spouse is having sexual relations with someone else. A paranoid disorder can seriously impair an individual's social or marital functioning, but the remainder of his thinking remains clear and orderly, his intellectual functioning is impaired only minimally or not at all, and the core of his personality remains intact. The treatment of persons with paranoid disorders involves the use of antipsychotic medications, frequently on a long-term-maintenance basis.

Mood disorders. These disorders are usually restricted to just two abnormalities of mood—depression and elation, or mania.

Depression is characterized by a sad or hopeless mood, pessimistic thinking, a loss of enjoyment and interest in one's usual activities and pastimes, reduced energy and vitality, increased fatigue, slowness of thought and action, change of appetite, and disturbed sleep. Depression must be distinguished from the grief and low spirits felt in reaction to the death of a loved one or some other unfortunate circumstance. The most dangerous consequence of severe depression is suicide.

Mania is characterized by an elated or euphoric mood, quickened thought and accelerated, loud, or voluble speech, overoptimism and heightened enthusiasm and confidence, inflated self-esteem, heightened motor activity, irritability, excitement, and a decreased need for sleep. The manic individual may injure himself, commit illegal acts, or suffer financial losses due to the poor judgment and risk-taking behaviour he displays when in the manic state.

Major mood disorders. The DSM-IV defines two major, or severe, mood disorders: bipolar disorder and major depression. A person with bipolar disorder, previously known as manic-depressive disorder, typically experiences discrete episodes of depression and then of mania lasting for a few weeks or months, with intervening periods of complete normality. The sequence of depression and mania can vary

widely from person to person and within one individual, with either mood abnormality predominating in duration and intensity. Depressive mood swings typically occur more often and last longer than manic ones, though there are persons who have episodes only of mania. Individuals with bipolar disorder frequently also show psychotic symptoms such as delusions, hallucinations, paranoia, or grossly bizarre behaviour.

The lifetime risk of developing bipolar disorder is about 1 percent and is about the same for men and women. The onset of the illness often occurs around the age of 30, and the illness persists over the long term. The predisposition to developing bipolar disorder is partly genetically inherited. Antipsychotic medications are used for the treatment of acute or psychotic mania. Mood-stabilizing agents such as lithium and several antiepileptic medications have proved effective in both treating and preventing recurrent attacks of mania.

Severe and long-lasting depression without the presence of mania is classified by the DSM-IV as major depressive disorder. Depression is a much more common illness than mania, and there are indeed many sufferers from depression who have never experienced mania. Major depressive disorder may occur as a single episode, or it may be recurrent. It may also exist with or without melancholia and with or without psychotic features. Melancholia implies the biological symptoms of depression: early-morning waking; daily variations of mood with depression most severe in the morning; loss of appetite and weight; constipation; and loss of interest in love and sex. Melancholia is a particular depressive syndrome that is relatively more responsive to somatic treatments such as medications (*e.g.*, Prozac, Paxil, and Zoloft) and electroconvulsive therapy (ECT).

It is estimated that women experience depression about twice as often as men. While the rates for major depression in men increase with age, the peak for women is between the ages of 35 and 45. There is a serious risk of suicide with the illness; of those who have a severe depressive disorder, about one-sixth eventually kill themselves. The loss of one's parents while a child or other childhood traumas or deprivations can increase a person's vulnerability to depression later in life, and stressful life events, especially where some type of loss is involved, are, in general, potent precipitating causes of the illness. It seems that both psychosocial and biochemical mechanisms are important in causing depression. Of the latter factor, the best-supported hypotheses suggest that the faulty regulation of the release of one or more neurotransmitters (*e.g.*, serotonin, dopamine, and norepinephrine) where the transmission of nerve impulses takes place is the basic cause, with a deficiency of the neurotransmitters resulting in depression and an excess causing mania. The treatment of major depressive episodes usually requires antidepressant medications. Electroconvulsive therapy may also be helpful, as may cognitive, behavioral, and interpersonal psychotherapies.

Additional mood disorders. A less severe manifestation of bipolar disorder, in which the mood swings are present but not as extreme, is termed cyclothymic disorder. The prevailing mood swings are established in adolescence and continue throughout adult life.

Dysthymic disorder, or depressive neurosis, may occur on its own, but it more commonly appears along with other neurotic symptoms such as anxiety, phobia, and hypochondriasis. Where there are clear external grounds for a person's unhappiness, a dysthymic disorder is considered to be present when the depressed mood is disproportionately severe or prolonged in regard to the distressing experience, when there is a preoccupation with the precipitating situation, when the depression continues even after removal of the provocation, and when it impairs the individual's ability to cope with the specific stress.

At any time, depressive symptoms may be present in one-sixth of the population, more commonly in women than men. Loss of self-esteem, feelings of helplessness and hopelessness, and loss of cherished possessions are also seen as important causes of minor depression. Psychotherapy is the treatment of choice for both dysthymic disorder and cyclothymic disorder, although antidepressant medications

Depression

or mood-stabilizing agents are often beneficial. Symptoms must be present for at least two years in order for a diagnosis of dysthymic or cyclothymic disorder to be made.

Anxiety disorders. Anxiety is defined as a feeling of fear, dread, or apprehension that arises without a clear or appropriate justification. Some experts differentiate anxiety from true fear in that the latter is experienced in response to an actual threat or danger, such as those to one's physical safety. Anxiety, on the other hand, may arise in response to apparently innocuous situations or may be out of proportion to the actual degree of the external stress. Anxiety also frequently arises as a result of subjective emotional conflicts of whose nature the person himself may be unaware. Generally, intense, persistent, or chronic anxiety that is not justified in response to real-life stresses and that interferes with the individual's functioning is regarded as a manifestation of mental disorder. Anxiety is a symptom of many mental disorders, including schizophrenia, obsessive-compulsive disorders, posttraumatic stress disorders, and so on, but, in phobias and other anxiety disorders proper, anxiety is the primary and frequently the only symptom.

The symptoms of anxiety are physical, psychological, and behavioral. Anxiety, especially during panic attacks, can manifest itself in a distinctive set of physical signs that arise from overactivity of the sympathetic nervous system or from tension in skeletal muscles. The sufferer experiences palpitations, dry mouth, dilation of the pupils, shortness of breath, sweating, abdominal symptoms, tightness in the throat, trembling, and dizziness. Aside from the actual feelings of dread and apprehension, the psychological symptoms include irritability, difficulty concentrating, and restlessness. Anxiety may also be manifested in avoidance behaviour (e.g., running away from the feared object or situation).

Phobic disorder. Phobias are neurotic states accompanied by intense dread of certain objects or situations that would not normally have such an effect. This type of anxiety is associated with a strong desire to avoid the dreaded object or situation. About 6 per 1,000 of the population suffer from a phobic disorder. There is a tendency for phobic symptoms, whatever their nature, to persist for many years unless treated, and the avoidance behaviour they produce can seriously limit the affected individual's movements and his social or occupational functioning.

People can have phobias about many different kinds of objects or situations, but three main divisions of phobic syndromes are made by the DSM-IV: specific phobia, agoraphobia, and social phobia. Individuals with specific phobias may intensely fear a specific object or situation—for example, cats or thunderstorms. They have anxious thoughts upon anticipating contact with an object or event—for instance, upon hearing the weather forecast—and they try to avoid the object, as in staying indoors in order not to encounter a cat. Typically, individuals with agoraphobia have an intense fear of being alone in or being unable to escape from a public place or some other setting outside the home, such as a crowded bus or a supermarket. A social phobia is present when individuals have extreme anxiety in a social situation where they are under the scrutiny of others, such as eating in a restaurant or speaking at a meeting.

The treatment of phobic disorders is best approached by the use of behavioral therapy. Dynamic psychotherapy and anti-anxiety medications may be effective in some cases.

Panic disorder and generalized anxiety disorder. Anxiety disorders in which the anxiety is not aroused by any specific object or situation can basically be subsumed under the headings of panic disorder and generalized anxiety disorder. Panic attacks are characterized by the sudden onset of intense or overwhelming anxiety accompanied by any of the aforementioned physical signs, such as difficulty in breathing, sweating, palpitations, and so on. The fear and apprehension experienced in such attacks sometimes mount to what are known as feelings of doom. Clear precipitating circumstances may produce the initial feelings of intense anxiety. The panic attack may last for about 15 minutes and often recurs, either infrequently or several times a week. The disorder usually starts in young adulthood and may persist for many years.

A diffuse and persistent feeling of anxiety associated with no particular object or situation is termed general, or free-floating, anxiety and is classified by the DSM-IV as generalized anxiety disorder. General anxiety is usually milder and less intense than in panic attacks, but it is longer-lasting and may persist for several months or years or on a recurrent basis. The most effective treatments vary according to the type of disorder and the individual. Psychotherapy and anti-anxiety medications are often useful in treating generalized anxiety disorder and panic disorder.

Obsessive-compulsive disorder. In this condition an individual experiences obsessions or compulsions or both. Obsessions are recurring words, thoughts, ideas, or images that, rather than being experienced as voluntarily produced, seem to invade a person's consciousness despite his attempts to ignore, control, or suppress them. The obsessional thought or topic is perceived by the individual as inappropriate or senseless; the idea is recognized as both alien to his nature and yet as coming from inside himself. An obsession can take the form of a recurrent and vivid fantasy that is often obscene, disgusting, or senseless. The person with obsessional ruminations holds endless debates over mundane matters inside his head (e.g., "Did I forget to lock the front door behind me?").

Obsessions in turn are frequently linked to compulsions. These are urges or impulses to perform repetitive acts that are apparently meaningless, unnecessary, stereotyped, or ritualistic. The compulsive person knows that the act to be performed is meaningless or unnecessary, but failure or refusal to perform it brings on a mounting tension or anxiety that is temporarily relieved once the act has been performed. Obsessional ruminations thus directly produce compulsive behaviour (e.g., repeatedly checking and relocking an already secure front door). Most compulsive acts have a simple, ritualistic character and can involve checking, touching, hand washing, or the repetition of particular words or phrases.

Medications, psychotherapy, and behavioral therapy are selectively successful in treating obsessive-compulsive disorders, depending on the individual. An older medication, clomipramine, and newer selective serotonin reuptake inhibitor (SSRI) medications have proved effective in reducing the symptoms in a large proportion of individuals with obsessive-compulsive disorder.

Posttraumatic stress disorder. In this condition symptoms develop in an individual after he has experienced a psychologically traumatic event. It is a category in the DSM-IV classification but is not different in its symptomatology from other anxiety conditions; the distinctive feature is the presence of external trauma. The traumatic events can include serious automobile accidents, rape or assault, military combat, torture, incarceration in a concentration or death camp, and such natural disasters as floods, fires, and earthquakes.

A feature of this condition is the person's reexperiencing of the traumatic event in nightmares and in intrusive daytime fantasies. Increased arousal is common, with insignificant events, such as a knock at the door, precipitating a sudden terrifying recollection and an exaggerated startle response. Other symptoms include emotional numbing, avoidance of stimuli associated with the trauma, a diminished ability to enjoy activities or relationships that were previously pleasurable, and difficulty in sleeping. Long-term symptoms of distress, marital and family problems, difficulties at work, and the abuse of alcohol and other drugs are characteristic impairments caused by the disorder.

The marked emotional symptoms may persist long after the actual occurrence of the traumatic event. Some persons are more liable than others to develop the disorder, depending on personality traits, previous psychological disturbances, age, and genetic predisposition. Psychotherapy is the basic approach used in treating this disorder.

Somatoform disorders. In these conditions the physical symptoms of the person suggest the presence of organic disease, but no such organic disorder can be found upon physical examination and investigation, and instead there is evidence that the symptoms are caused by psychological factors. The production of these symptoms is not under

Physical
pain
caused by
psychologi-
cal conflict

voluntary control. The terms *hypochondriasis* and *hysteria* were traditionally used to designate these disorders.

Somatization disorder. This disorder was previously designated Briquet's syndrome; its essential features consist of multiple recurrent physical complaints made over many years and starting in young adult life or adolescence. The individual demands medical attention, but no organic cause is found. The symptoms invariably occur in many different bodily systems—for instance, back pains, dizziness, indigestion, difficulty with vision, and partial paralysis—and the symptoms may follow trends in health concerns among the public.

The condition is relatively common and occurs in about 1 percent of adult women. Males rarely exhibit this disorder. There are no clear etiological factors. Treatment involves not agreeing with the person's inclination to attribute organic causes to the symptoms and ensuring that physicians and surgeons do not cooperate with the person in seeking excessive diagnostic procedures or surgical remedies for the complaints.

Conversion disorder. This disorder was previously labeled *hysteria*. Its symptoms are a loss of or alteration in physical functioning, which may include paralysis. The physical symptoms occur in the absence of organic pathology and are instead apparently the expression of an underlying emotional conflict. The characteristic motor symptoms of conversion disorder include the paralysis of the voluntary muscles of an arm or leg, tremor, tics, and other disorders of movement or gait. The neurological symptoms may be widely distributed and may not conform with medical knowledge of nerve distribution. Blindness, deafness, loss of sensation in arms or legs, the feeling of "pins and needles," and an increased sensitivity to pain in a limb may also be present.

Symptoms usually occur in a setting of extreme psychological stress and appear suddenly. The course of the disorder is variable, with recovery often occurring in a few days but with symptoms persisting for years or decades in chronic cases that remain untreated.

The causation of conversion disorder has been linked with fixations (*i.e.*, arrested stages in the individual's early psychosexual development). Freud's theory that threatening or emotionally charged thoughts are repressed out of consciousness and converted into physical symptoms is still widely accepted. The treatment of conversion disorder thus requires psychological rather than pharmacological methods, notably the exploration of the individual's underlying emotional conflicts. Conversion disorder can also be considered as a form of "illness behaviour"; *i.e.*, the person uses the symptoms to gain a psychological advantage in social relationships, either by gathering sympathy or by being relieved of burdensome or stressful obligations and withdrawing from emotionally disturbing or threatening situations. Thus it may be advantageous to the person, in a psychological sense, to have the consequences of the symptoms.

Hypochondriasis. Hypochondriasis is a preoccupation with physical signs or symptoms that the person unrealistically interprets as abnormal, leading to the fear or belief that he is seriously ill. There may be fears about the development of physical or mental symptoms without any such existing, a belief that actual but minor symptoms are of dire consequence, or an experience of normal bodily sensations as threatening symptoms. A thorough physical examination may find no organic cause for the physical signs the individual is concerned about, but the examination fails to relieve his unrealistic fears about having a serious disease. The symptoms of hypochondriasis may occur with mental illnesses other than anxiety, such as depression or schizophrenia.

The onset of this disorder may be associated with precipitating factors such as an actual organic disease with physical and psychological aftereffects—*e.g.*, coronary thrombosis in a previously fit man. It often begins during the fourth and fifth decades of life but is also common at other times, during pregnancy, for example. Treatment aims to provide understanding and support and to reinforce healthy behaviour; antidepressant medications may be used to relieve depressive symptoms.

Psychogenic pain disorder. In psychogenic pain disorder the main feature is a persistent complaint of pain in the absence of organic disease and with evidence of a psychological cause. The pattern of pain may not conform to the known anatomic distribution of the nervous system. Psychogenic pain may occur as part of hypochondriasis or as a symptom of a depressive disorder. Appropriate treatment depends on the context of the symptom.

Dissociative disorders. A dissociative disorder is a syndrome in which one or a group of mental processes are split off, or dissociated, from the rest of the psychological apparatus so that their function is lost, altered, or impaired. Dissociative symptoms have often been regarded as the mental counterparts of the physical symptoms displayed in conversion disorders. Since the dissociation may be an unconscious mental attempt to protect the individual from threatening impulses or emotions that are repressed, the conversion into physical symptoms and the dissociation of mental processes can be seen as related defense mechanisms arising in response to emotional conflict.

In dissociative disorders there is a sudden, temporary alteration in the person's consciousness, sense of identity, or motor behaviour. There may be an apparent loss of memory of previous activities or important personal events, with amnesia for the episode itself after recovery. These are rare conditions, and it is important to exclude organic causes.

In dissociative amnesia there is a sudden loss of memory that may appear total; the individual can remember nothing about his previous life or even his name. The amnesia may be localized to a short period of time associated with a traumatic event, or it may be selective, affecting the person's recall of some, but not all, of the events during a particular time. In psychogenic fugue, the individual wanders away from home or place of work and assumes a new identity; he cannot remember his previous identity and upon recovering cannot recall the events that occurred while he was in the fugue state. In many cases the disturbance lasts only a few hours or days and involves only limited travel. Severe stress frequently triggers this disorder.

Dissociative identity disorder, previously called multiple personality disorder, is a rare and remarkable dissociative disorder in which two or more distinct and independent personalities develop in a single individual. Each of these personalities inhabits the person's conscious awareness to the exclusion of the others at particular times. This disorder frequently arises as a result of traumas suffered during childhood and is best treated by psychotherapy, which seeks to reunite the various personalities into a single, integrated one.

In depersonalization disorder a person feels or perceives his body or self as being unreal, strange, altered in quality, or distant. This state of self-estrangement may take the form of feeling as if one is machine-like, is living in a dream, or is not in control of one's actions. Derealization, or feelings of unreality concerning objects outside one's self, often occurs at the same time. Depersonalization disorder may occur alone in neurotic persons but is more often associated with phobic, anxiety, or depressive symptoms. It most commonly occurs in younger women and may persist for many years. Persons find the experience of depersonalization intensely difficult to describe and often fear that others will think them insane. Organic conditions, especially temporal lobe epilepsy, must be excluded before making a diagnosis of neurosis when depersonalization occurs. As with other neurotic syndromes, it is more common to see many different symptoms than to see depersonalization alone.

The causes of depersonalization are obscure, and there is no specific treatment for it. When the symptom arises in the context of another psychiatric condition, treatment is aimed at that illness.

Personality disorders. Personality is the characteristic way in which an individual thinks, feels, and behaves; it accounts for the ingrained behaviour patterns of the individual and allows the prediction of how he will act in particular circumstances. Personality embraces a person's moods, attitudes, and opinions and is most clearly expressed in interactions with other people. A personality disorder is a deeply ingrained, long-enduring, maladaptive,

Multiple
personality
disorder

and inflexible pattern of thinking, feeling, and behaving that either significantly impairs an individual's social or occupational functioning or causes subjective distress.

A personality disorder may occur with another psychiatric condition or on its own. The causes of personality disorders are obscure. There is undoubtedly a constitutional and therefore hereditary element in determining personality type. Psychological and environmental factors are also important in causation.

Some generally accepted types of personality disorder are listed below. It is important to recognize that simply exhibiting the trait or even having it to an abnormal extent is not enough to constitute disorder—for that, the degree of abnormality must cause disturbance to the individual or to society. Because personality traits are, by definition, virtually permanent, these disorders are only partially, if at all, amenable to treatment. The most effective treatment combines various types of group, behavioral, and cognitive psychotherapy. The behavioral manifestations of personality disorders often tend to diminish in their intensity in middle and old age.

Paranoid personality disorder. In this disorder there is a pervasive and unjustified suspiciousness and mistrust of others, whose words and actions are misinterpreted as having special significance for, and as being directed against, the individual. Sometimes such people are guarded, secretive, aggressive, quarrelsome, and litigious, and they are excessively sensitive to the implied criticism of others.

Schizoid personality disorder. In this disorder there is a disinclination to interact with others; the individual appears aloof, withdrawn, indifferent, unresponsive, and uninterested. Such a person prefers solitary to gregarious pursuits, involvement with things rather than with people, and often appears humourless or dull.

Schizotypal personality disorder. This category has been used to describe people who show various oddities or eccentricities of thought, speech, perception, or behaviour (e.g., bizarre fantasies or persecutory delusions) but whose symptoms are not severe enough to be labeled as schizophrenic.

Obsessive-compulsive personality disorder. A person with this disorder shows prominent overscrupulous, perfectionistic traits that are expressed in feelings of insecurity, self-doubt, meticulous conscientiousness, indecisiveness, and rigidity of behaviour. The person is preoccupied with rules, procedures, and efficiency, is overly devoted to work and productivity, and is usually deficient in the ability to express warm or tender emotions. This disorder is more common in men and is in many ways the antithesis of antisocial personality disorder (see below).

Avoidant personality disorder. A person with this disorder shows extreme sensitivity to rejection and may lead a socially withdrawn life. Nevertheless, the person is not asocial and does show a great desire for companionship but needs unusually strong guarantees of uncritical acceptance. Persons with this disorder are commonly described as having an "inferiority complex."

Dependent personality disorder. This category is used to describe individuals who subordinate their own needs and responsibility over major areas of their lives to the control of others. Persons with this disorder lack self-confidence and may experience intense discomfort when alone. Women are affected with this condition more often than men.

Histrionic personality disorder. People with this disorder are overly dramatic and intensely expressive, egocentric, highly reactive, and excitable. They exhibit a tendency to call attention to themselves, crave novelty and excitement, and may be dependent and suggestible. This disorder is more common in women than men.

Antisocial personality disorder. This disorder is marked by a personal history of chronic and continuous antisocial behaviour, in which the rights of others are violated, and by poor or nonexistent job performance. It is manifested in persistent criminality, sexual promiscuity or aggressive sexual behaviour, and drug use. People with this disorder are impulsive, mendacious, irresponsible, and callous; they feel no guilt over their antisocial acts and fail to learn from their mistakes. The symptoms are usually prominent by

adolescence. In order to be diagnosed with antisocial personality disorder, an individual must exhibit such symptoms before the age of 15. Antisocial personalities are less liable to criminal acts as they grow older, but there remains a high risk of suicide, accidental death, drug or alcohol abuse, and a tendency toward interpersonal problems. Antisocial personality disorder is more common in men than in women.

Narcissistic personality disorder. A person with this disorder has a grandiose sense of self-importance and a preoccupation with fantasies of success, power, and achievement.

Borderline personality disorder. Borderline personality disorder is characterized by an extraordinarily unstable mood and self-image. Individuals with this disorder may exhibit intense episodes of anger, depression, or anxiety. Self-injury is also a common symptom of this disorder.

Psychosexual disorders. The following section discusses disorders of gender identity and preferences for unusual or bizarre sexual practices or objects. Psychosexual dysfunctions such as impotence are treated in the article *SEXUAL BEHAVIOUR, HUMAN*.

Gender identity disorder. In gender identity disorder a person feels a discrepancy between his anatomic sex and the gender that he ascribes to himself. This disorder is much more common in males than females. The individual claims that he is a member of the opposite sex—"a female mind trapped in a male body." An individual with gender identity disorder may assume the dress and behaviour of and participate in activities commonly associated with the opposite sex and may even use hormones and surgery to achieve the physical characteristics of the opposite sex. The cause of the condition is unknown. Individuals with this disorder have a significant risk of developing depression and an increased risk of suicide. Psychiatric treatment is generally supportive in type. Persons with gender identity disorder may choose to have sex reassignment surgery, a procedure in which the body, including the genitals, is surgically altered to look like that of the opposite sex. Sex-reassignment surgery requires following carefully prescribed standards.

Paraphilias. Paraphilias, or sexual deviations, are defined as unusual fantasies, urges, or behaviours that are recurrent and sexually arousing. These urges must occur for at least six months and cause distress to the individual in order to be classified as a paraphilia. In fetishism inanimate objects (e.g., shoes) are the person's sexual preference and means of sexual arousal. In transvestism the recurrent wearing of clothes of the opposite sex is performed to achieve sexual excitement. In pedophilia an adult has sexual fantasies about or engages in sexual acts with a prepubertal child of the same or opposite sex.

In exhibitionism repeated exposure of the genitals to an unsuspecting stranger is used to achieve sexual excitement. In voyeurism observing the sexual activity of others repeatedly is the preferred means of sexual arousal. In sadomasochism the individual achieves sexual excitement as either the recipient or the provider of pain, humiliation, or bondage.

There are, of course, other unusual sexual objects or acts that may be used for gratification. The causes of these conditions are generally not known. Behavioral, psychodynamic, and pharmacological methods have been used with varying efficacy to treat these disorders.

Disorders usually first evident in infancy, childhood, or adolescence. Children are usually referred to a psychiatrist or therapist because of a parent's or some other adult's complaints or concern over the child's behaviour or development. Family problems, particularly difficulties in the parent-child relationship, are often an important causative factor in the symptomatic behaviour of the child. For a child psychiatrist, the observation of behaviour is especially important, as the child may not be able to express his feelings in words. Isolated psychological symptoms are extremely common in children. Boys are affected twice as often as girls.

Attention-deficit disorders. Children with these disorders show a degree of inattention and impulsiveness that is markedly inappropriate for their stage of development.

Gross overactivity in children can have many causes, including anxiety, conduct disorder (see below), or the effects of living in institutions. Learning difficulties and antisocial behaviour may occur secondarily. This syndrome is 10 times more common in boys than in girls.

Conduct disorders. These are the most common psychiatric disorders in older children and adolescents, accounting for nearly two-thirds of disorders in those aged 10 and 11 years. Abnormal conduct more serious than ordinary childish mischief persistently occurs; lying, disobedience, aggression, truancy, delinquency, and deterioration of work may occur at home or at school. Vandalism, drug and alcohol abuse, and early sexual promiscuity may also occur. The most important causative factors are the family background; broken homes, unstable and rejecting families, institutional care in childhood, and a poor social environment are frequently present in such cases.

Anxiety disorders. Neurotic or emotional disorders in children are similar to the adult conditions except that they are often less clearly differentiated. In anxiety disorders of childhood, the child is fearful, timid with other children, and overdependent and clinging toward the parents. Physical symptoms, sleep disturbance, and nightmares occur. Separation from the parent or from the home environment is a major cause of this anxiety.

Eating disorders. Anorexia nervosa usually starts in late adolescence and is about 20 times more common in girls than boys. This disorder is characterized by a failure to maintain one's body weight at the normal level for an individual's age and height; weight loss is at least 15 percent of the ideal body weight. Weight loss occurs because of an intense desire to be thin, a fear of gaining weight, or a disturbance in the way in which one sees the body weight or shape. Postmenarchal females with anorexia nervosa usually experience amenorrhea (*i.e.*, the absence of at least three consecutive menstrual periods). Though grossly thin, individuals with anorexia nervosa nevertheless believe that they are fat. They will go to enormous lengths to resist eating food and to lose weight, including food avoidance, self-induced vomiting, and vigorous exercise. Medical complications of anorexia nervosa can be life-threatening.

The condition appears to start with an individual's voluntary control of food intake in response to social pressures such as peer conformity. The disorder is exacerbated by troubled relations within the family. It is much more common in developed, wealthy societies and in girls of higher socioeconomic class. Treatment includes persuading the person to accept and cooperate with medical therapy, achieving weight gain, and helping the person maintain weight by psychological and social therapy.

Bulimia nervosa is characterized by excessive overeating binges combined with inappropriate methods of stopping weight gain such as self-induced vomiting, laxatives, or diuretics. Repetition of the cycle can lead to serious medical complications, such as dental decay or dehydration, and can be fatal. Like persons with anorexia nervosa, individuals who have bulimia nervosa show great concern for their weight and body shape; the majority, however, are close to their proper weight. Anorexia nervosa and bulimia nervosa may occur simultaneously in the same person. Depression, anxiety, and low self-esteem may be associated with the disorder.

Autistic disorder. Psychotic disorders are very rare in childhood, and of these about one-half are cases of autistic disorder. Boys are affected three times as often as girls. Autistic disorder begins in the first two years of life and is more common in the upper socioeconomic classes. The child exhibits an inability to make warm emotional relationships, has a severe speech and language disorder, and exhibits a desire for routine in which he shows distress if thwarted from his stereotyped behaviour. There is some evidence to support genetic and organic factors in the causation of autistic disorder. Treatment involves management of the abnormal behaviour, training in life skills and occupational activities, and counseling for the family.

Other childhood disorders. Stereotyped movement disorders involve the exhibition of tics in differing patterns. A tic is an involuntary, purposeless jerking movement of a group of muscles or the involuntary production of noises

or words. Tics may affect the face, head, and neck or, less commonly, the limbs or trunk. Tourette's syndrome is typified by multiple tics and involuntary vocalization, which sometimes includes the uttering of obscenities.

Other physical symptoms that are often listed among psychiatric disorders of childhood include stuttering, enuresis (*i.e.*, the repeated involuntary emptying of urine from the bladder during the day or night), encopresis (*i.e.*, the repeated voiding of feces into inappropriate places), sleepwalking, and night terror. These symptoms are not necessarily evidence of emotional disturbance or of some other mental illness. Behavioral methods of treatment are usually effective.

Other mental disorders. **Factitious disorders.** Factitious disorders are characterized by physical or psychological symptoms that are voluntarily self-induced; they are distinguished from conversion disorder, in which the physical symptoms are produced unconsciously. In factitious disorders, although the person's attempts to create or exacerbate the symptoms of an illness are voluntary, such behaviour is neurotic in that the individual is unable to refrain from it; *i.e.*, his goals, whatever they may be, are involuntarily adopted. In malingering, by contrast, the person stimulates or exaggerates an illness or disability to obtain some kind of discernible personal gain or to avoid an unpleasant situation; *e.g.*, a prison inmate may simulate madness to obtain more comfortable living conditions. It is important to recognize factitious disorders as evidence of psychological disturbance.

Impulse-control disorders. Persons with these conditions demonstrate a failure to resist desires, impulses, or temptations to perform an act that is harmful to the individual or to others. The individual experiences a feeling of tension before committing the act and a feeling of release or gratification upon completing it. The behaviours involved include pathological gambling, pathological setting of fires (pyromania), pathological stealing (kleptomania), and recurrent pulling of hair (trichotillomania).

Adjustment disorders. These are conditions in which there is an inappropriate reaction to an external stress occurring within three months of the stress. The symptoms may be out of proportion to the degree of stress, or they may be maladaptive in the sense that they prevent the individual from coping adequately in his social or occupational setting. These disorders are often associated with other mood or anxiety disorders. (A.C.P.S./L.B.An.)

Treatment of mental disorders

HISTORICAL OVERVIEW

Early history. References to mental disorders in early Egyptian, Indian, Greek, and Roman writings show that the physicians and philosophers who contemplated problems of human behaviour regarded mental illnesses as a reflection of the displeasure of the gods or as evidence of demoniac possession. Only a few realized that individuals with mental illnesses should be treated humanely rather than exorcised, punished, or banished. Certain Greek medical writers, however, notably Hippocrates (flourished 400 BC), regarded mental disorders as diseases to be understood in terms of disturbed physiology. Hippocrates and his followers emphasized natural causes, clinical observation, and brain pathology in the study of mental disorders. Later Greek medical writers, including those who practiced in Imperial Rome, prescribed treatments for mental illness, including a quiet environment, work, and the use of drugs such as the purgative hellebore. It is probable that most people with psychoses during ancient times were cared for by their families and that those who were thought to be dangerous to themselves or others were detained at home by relatives or hired keepers.

During the early Middle Ages in Europe, primitive thinking about mental illness reemerged, and witchcraft and demonology were used to account for the symptoms and behaviour of people with psychoses. At least some of those who were deemed insane were looked after by the religious orders, which offered care for the sick generally. The empirical and quasi-scientific Greek tradition in medicine was maintained not by the Europeans but by the Muslim

Arabs, who are usually credited with the establishment of asylums for the mentally ill in the Middle East as early as the 8th century. In medieval Europe in general it seems that the mentally ill were allowed their freedom, provided they were not regarded as dangerous. The founding of the first hospital in Europe devoted entirely to the care of the mentally ill probably occurred in Valencia, Spain, in 1407–09, though this has also been said of a hospital established in Granada in 1366–67.

Brutal
treatment
of the
mentally ill

From the 17th century onward in Europe, there was a growing tendency to isolate deviant people, including the mentally ill, from the rest of society. Thus, such socially unwanted people as the mentally ill were confined together with the disabled, vagrants, and delinquents. Those of the mentally ill who were regarded as violent were often chained to the walls of prisons and were treated in a barbarous and inhumane way.

In the 17th and 18th centuries, the development of European medicine and the rise of empirical methods of medicoscientific inquiry were paralleled by an improvement in public attitudes toward the mentally ill. By the end of the 18th century, concern over the care of the mentally ill had become so great among educated people in Europe and North America that governments were forced to act. After the French Revolution the physician Philippe Pinel was placed in charge of the Bicêtre, the hospital for the mentally ill in Paris. Under Pinel's supervision a completely new approach to the care of mental patients was introduced. Chains and shackles were removed from the patients, and, in place of dungeons, the mentally ill were provided with sunny rooms and were permitted to exercise on the hospital grounds. Among other reformers were the British Quaker layman William Tuke, who established the York Retreat for the humane care of the mentally ill in 1796, and the physician Vincenzo Chiarugi, who published a humanitarian regime for his hospital in Florence in 1788. In the mid-19th century Dorothea Dix led a campaign to increase public awareness of the inhumane conditions that prevailed in American mental hospitals, and her efforts led to widespread reforms both in the United States and elsewhere.

The mental hospital era. Many hospitals for the mentally ill were built in the latter half of the 18th century. Some of them, like the York Retreat in England, were run on humane and enlightened lines, while others, like the York Asylum, gave rise to great scandal because of their brutal methods and filthy living conditions. In the mid-19th century an extensive program of mental hospital building was carried out in North America, Britain, and many of the countries of continental Europe. The hospitals housed poor mental patients, and their aim was to care for these individuals humanely and to relieve their families of the burden of caring for them. The approach was that of moral treatment, including work, the avoidance of physical methods of restraint, and respect for the individual patient. A widespread belief in the curability of mental illness at this time was a principal motivating factor behind such reform.

The mental hospital era was an age of reform, and there is no doubt that patients were treated much more humanely. The era produced a large number of segregated institutions in which a much higher proportion of the mentally ill was confined than previously. But the medical reformers' early hopes of successful cures were not vindicated, and by the end of the 19th century the hospitals had become overcrowded, and custodial care had replaced moral treatment.

The biological movement. Along with humanitarian reforms in hospital practice and treatment methods during the late 18th and 19th centuries, there was a resurgence of medical and scientific interest in psychiatric theory and practice. Fundamental strides were made during this period in establishing a scientific basis for the study of mental disorders. A long series of observations by clinicians in France, Germany, and England culminated in 1883 in a comprehensive classification of mental disorders by the German psychiatrist Emil Kraepelin. His classification system served as the basis for all subsequent ones, and the cardinal distinction he made between schizophrenia and bipolar disorder still stands.

Rapid advances in various branches of medicine led in the later 19th century to the expectation of discovering specific brain lesions that were thought to cause the various forms of mental disorder. While this research did not attain the results that were expected, the scientific emphasis was productive in that it did elucidate the gross and microscopic pathology of many brain disorders that can produce psychiatric disabilities. Nevertheless, many of the psychotic disorders, notably schizophrenia and bipolar disorder, frustrated the effort to find causative agents in cellular pathology. It became apparent that other explanations had to be found for the many puzzling aspects of mental disorders in general, and these explanations emerged in a wave of psychological rather than physical explanations.

Development of psychotherapy. Foremost among these was that of psychoanalysis, which originated in the work of the Viennese neurologist Sigmund Freud. Having studied under the French neurologist Jean-Martin Charcot, Freud originally used well-known techniques of hypnosis to treat patients suffering from what was then called hysterical paralysis and other neurotic syndromes. Freud and his colleague Josef Breuer observed that their patients tended to relive earlier life experiences that could be associated with the symptomatic expression of their illnesses. When these memories and the emotions associated with them were brought to consciousness during the hypnotic state, the patients showed improvement. Observing that most of his patients proved able to talk about such memories without being under hypnosis, Freud evolved the technique of free association (*i.e.*, production by the patients, aloud and without suppression or self-censorship of any kind, of the thoughts and feelings about whatever was uppermost in their minds) as a means of access to the unconscious. From this beginning Freud gradually developed what became known as psychoanalysis. Other features of the new procedure included the study of dreams, the interpretation of "resistances" on the part of the patient, and the handling by the therapist of transference (*i.e.*, the patient's feelings toward the analyst that reflected previously experienced feelings toward parents and other important figures in his early life). Freud's work, though complex and controversial in many of its aspects, laid the basis for modern psychotherapy in its use of free association and its emphasis on unconscious and irrational mental processes as causative factors in mental illness. This emphasis on purely psychological factors as a basis for both causation and treatment was to become the cornerstone of most subsequent psychotherapies.

Freud's
discoveries

Variations of the original psychoanalytic technique were introduced by several of Freud's colleagues who parted company with him. Analytic psychology, devised by Carl Jung, placed less emphasis on free association and more on the interpretation of dreams and fantasies. Special importance was given to the collective unconscious, a reservoir of shared unconscious wisdom and ancestral experience that entered consciousness only in symbolic form to influence thought and behaviour. Jungian analysts sought clues to their patients' problems in the archetypal nature of myths, stories, and dreams. Individual psychology, devised by Alfred Adler, emphasized the importance of the individual's drive toward power and of his unconscious feelings of inferiority. The therapist was concerned with the patient's compensations for his inferiority, as well as with his social relationships.

Development of physical and pharmacological treatments. During the early decades of the 20th century, the principal approaches to the treatment of mental disorders were psychoanalytically derived psychotherapies, used to treat people with neuroses, and custodial care in mental hospitals, for those with psychoses. But beginning in the 1930s these methods began to be supplemented by physical approaches using drugs, electroconvulsive therapy, and surgery. The first successful physical treatment in psychiatry was the induction of malaria in patients with a fatal form of neurosyphilis called general paresis. The malarial treatment stemmed from the observation that some psychotic patients improved during febrile illnesses. In 1933 the Polish psychiatrist Manfred Sakel reported that psychotic symptoms of patients with schizophrenia were improved by re-

Lobotomy

peated insulin-induced comas. (Neither of these treatments is in use today.) The treatment of symptoms of schizophrenia by convulsions, originally induced by the injection of camphor, was reported in 1935 by the psychiatrist Ladislaus Joseph von Meduna in Budapest. An improvement in this approach was the induction of convulsions by the passage of an electrical current through the brain, a technique introduced by Italian psychiatrists Ugo Cerletti and Lucio Bini in 1938. Electroconvulsive treatment was more successful in alleviating states of severe depression than in treating symptoms of schizophrenia. Psychosurgery, or surgery performed to treat mental illness, was introduced by the Portuguese neurologist António Egas Moniz in the 1930s. The procedure Moniz originated, leucotomy, or lobotomy, was widely performed during the next two decades in the treatment of patients with schizophrenia, intractable depression, and severe obsessional states. The procedure was later abandoned, however, largely because its therapeutic effects could be better obtained by the use of newly developed medications.

The decades after World War II were marked by the first safe and effective applications of medications in the treatment of mental disorders. Prior to the 1950s such sedative compounds as bromides and barbiturates had been used to quiet or sedate patients, but these drugs were general in their effect and did not target the specific symptoms of mood disturbances or psychotic disorders. Many of the medications that subsequently proved effective in treating such conditions were recognized serendipitously—*i.e.*, when researchers administered them to patients just to see what would happen, or when they were administered to treat one medical condition and were instead found to be helpful in alleviating the symptoms of a mental disorder.

The first effective pharmacological treatment of psychosis was the treatment of mania with lithium, introduced by the Australian psychiatrist J.F.J. Cade in 1949. Lithium, however, generated little interest until its dramatic effectiveness in the maintenance treatment of bipolar affective disorder was reported in the mid-1960s. Chlorpromazine, the first of a long series of highly successful antipsychotic drugs, was synthesized in France in 1950 during work on antihistamines. It was used in anesthesia before its antipsychotic and tranquilizing effects were reported in France in 1952. The first tricyclic antidepressant drug, imipramine, was originally designed as an antipsychotic drug and was investigated by the Swiss psychiatrist Roland Kuhn. He found it ineffective in treating symptoms of schizophrenia but observed its antidepressant effect, which he reported in 1957. A drug used in the treatment of tuberculosis, iproniazid, was found to be effective as an antidepressant in the mid-1950s. It was the first monoamine oxidase inhibitor to be used in psychiatry. The first modern anxiety-relieving drug was meprobamate, which was originally introduced as a muscle relaxant. It was soon overtaken by the pharmacologically rather similar but clinically more effective chlordiazepoxide, which was synthesized in 1957 and marketed as Librium in 1960. This drug was the first of the extensively used benzodiazepines. These and other drugs had a revolutionary impact not only on psychiatry's ability to relieve the symptoms and suffering of people with a wide range of mental disorders but also on the institutional care of the mentally ill.

Deinstitutionalization. Between about 1850 and 1950 there was a steady increase in the number of patients staying in mental hospitals. In England and Wales, for example, there were just over 7,000 such patients in 1850, nearly 120,000 in 1930, and nearly 150,000 in 1954. Then a steady decline in the number of occupied beds began, reaching just over 100,000 in 1970 and 75,000 in 1980, a decrease of almost 50 percent. The same process began in the United States in 1955 but continued at a more rapid rate. The decrease, from just under 560,000 in 1955 to just over 130,000 in 1980, was one of more than 75 percent. In both countries it became official policy to replace mental hospital treatment with community care, involving district general hospital psychiatric units in Britain and local mental health centres in the United States. This dramatic change can be partly attributed to the introduction of antipsychotic medications, which drastically changed the

atmosphere of mental hospital wards. With the recovery of lucidity and calmness, many patients with psychoses could return to their homes and live at least a partially normal existence instead of spending their lives sequestered in mental hospitals. The wholesale release of mental patients into the community was not without problems, however, since many areas lacked the facilities to support and maintain such patients, many of whom thus received inadequate care.

Development of behaviour therapy. In the 1950s and '60s a new type of therapy, called behaviour therapy, was developed. In contrast to the existing psychotherapies, its techniques were based on theories of learning derived from research on classical conditioning carried out by Ivan Pavlov and others and from the work of such American behaviourists as John B. Watson and B.F. Skinner. Behavioural therapy developed when the theoretical principles that were originally developed from experiments with animals were applied to the treatment of patients.

In 1920 Watson experimentally induced a phobia of rats in a small boy, and in 1924 Mary Cover Jones reported the extinction of phobias in children by gradual desensitization. Modern behaviour therapy began with the description by the South African psychiatrist Joseph Wolpe of his technique of systematically desensitizing a patient with phobias, beginning by exposing him to the least and gradually progressing to the most feared object or situation. Behavioral therapies were more quickly adopted in Europe than in the United States, where psychoanalytic precepts had exercised a particular dominance over psychiatry, but by the 1980s behavioral therapies were also well established in the United States.

The mental health profession in the late 20th century. There has been a great increase in the number of mental health professionals since World War II. In the United States the number of psychiatrists was 3,000 in 1939 but had increased to more than 50,000 by the early 1990s. Nonmedical mental health professionals have also increased substantially in number. Clinical psychologists, who at one time largely administered psychometric tests, now also provide psychotherapy and behaviour therapy. Psychiatric social workers also have become psychotherapists and play prominent roles in mental health centres. There are new roles for nurses, including behaviour therapy and the management of chronic mental illness in the community.

Psychotherapy retains a major role in the mental health profession. Subsequent to the development of psychoanalysis, the varieties of psychotherapy have increased and multiplied to more than 250 different therapies. The repertoire of medications used in the treatment of mental illness has continued to grow as new drugs are developed or new applications of existing ones are discovered. Research on the biochemical and genetic causes of mental disease continues to make headway in explicating the causes of various disorders. The triad of psychotherapy, medications, and counseling affords an unprecedented array of approaches, techniques, and procedures for alleviating the symptoms of people with mental disorders.

PHYSIOLOGICAL TREATMENTS

Pharmacological treatments. *Antipsychotic agents.* Antipsychotic medications, which are also known as neuroleptics and major tranquilizers, belong to several different chemical groups but are similar in their therapeutic effects. These medications have a calming effect that is valuable in the relief of agitation, excitement, and violent behaviour in persons with psychoses. The drugs are quite successful in reducing the symptoms of schizophrenia, mania, and delirium, and they are used in combination with antidepressants to treat psychotic depression. The drugs suppress hallucinations and delusions, alleviate disordered or disorganized thinking, improve the patient's lucidity, and generally make an individual more receptive to psychotherapy. Patients who have previously been agitated, intractable, or grossly delusional become noticeably calmer, quieter, and more rational when maintained on these drugs. The medications have enabled many patients with episodic psychoses to have shorter stays in hospitals,

and many other patients who would have been permanently confined to institutions are able to live in the outside world because of the drugs. The antipsychotics differ in their unwanted effects: some are more likely to make the patient drowsy, some to alter blood pressure or heart rate, and some to cause tremor or slowness of movement. Some of the most widely used antipsychotic drugs are shown in Table 1.

Table 1: Examples of Antipsychotic Drugs

drug	representative trade name
chlorpromazine	Thorazine
trifluoperazine	Stelazine
thiothixine	Navane
haloperidol	Haldol
clozapine	Clozaril
olanzepine	Zyprexa
risperidone	Risperdal
quetiapine	Seroquel

Treatment of schizophrenia

In the treatment of schizophrenia, antipsychotic drugs partially or completely control such symptoms as delusions and hallucinations. They also protect the patient who has recovered from an acute episode of the mental illness from suffering a relapse. The newer antipsychotic medications also treat social withdrawal, apathy, blunted emotional capacity, and the other psychological deficits characteristic of the chronic stage of the illness.

No single drug seems to be outstanding in the treatment of schizophrenia. In an individual patient, one drug may be preferred to another because it produces less severe unwanted effects, and the dose of any one drug needed to produce a therapeutic effect varies widely from patient to patient. Because of these individual differences, it is common for psychiatrists to substitute a drug of a different chemical group when one drug has been shown to be ineffective despite its use in adequate dosage for several weeks.

In an acute psychotic episode a drug such as chlorpromazine, olanzepine, or haloperidol usually has a calming effect within a day or two. The control of psychotic symptoms such as hallucinations or disordered thinking may take weeks. The appropriate dosage has to be determined for each patient by cautiously increasing the dose until a therapeutic effect is achieved without unacceptable side effects.

It is not known exactly how antipsychotic medications work. One theory is that they block dopamine receptors in the brain. Dopamine is a neurotransmitter—*i.e.*, a chemical messenger produced by certain nerve cells that influence the function of other nerve cells by interacting with receptors in their cell membranes. Since schizophrenia may be caused by either the excessive release of or an increased sensitivity to dopamine in the brain, the effects of antipsychotic drugs may be due to their ability to block or inhibit dopamine transmission.

Dopamine receptor blockade is certainly responsible for the main side effects of first-generation antipsychotic medications. These symptoms, which are termed extrapyramidal symptoms (EPS), resemble those of Parkinson's disease and include tremor (shakiness) of the limbs, bradykinesia (slowness of movement with loss of facial expression, absence of arm-swinging during walking, and a general muscular rigidity), dystonia (sudden, sustained contraction of muscle groups causing abnormal postures), akathisia (a subjective feeling of restlessness leading to an inability to keep still), and tardive dyskinesia (involuntary movements, particularly involving the lips and tongue). Most extrapyramidal symptoms disappear when the drug is withdrawn. Tardive dyskinesia occurs late in the drug treatment and in about half of the cases persists even after the drug is no longer used. There is no satisfactory treatment for severe tardive dyskinesia.

Antianxiety agents. The drugs most commonly used in the treatment of anxiety are the benzodiazepines, which have replaced the barbiturates because of their vastly greater safety. One advantage of benzodiazepines is that they are less dangerous in overdose than barbiturates. A large number of benzodiazepines have been marketed, and

Table 2: Commonly Used Benzodiazepines in the Treatment of Anxiety

drug	representative trade name
lorazepam	Ativan
oxazepam	Serax
alprazolam	Xanax
chlordiazepoxide	Librium
diazepam	Valium
chlorazepate	Tranxene

the more common ones are listed in Table 2. Benzodiazepines differ from each other in duration of action rather than effectiveness. Smaller doses have a calming effect and alleviate both the physical and the psychological symptoms of anxiety. Larger doses induce sleep, and some benzodiazepines are marketed as hypnotics. The benzodiazepines have become among the most widely prescribed drugs in the developed world, and controversy has arisen over their excessive use by the public.

The side effects of these medications are usually few—most often drowsiness and unsteadiness. Benzodiazepines are not lethal even in very large overdoses, but they increase the sedative effects of alcohol and other drugs. The benzodiazepines are basically intended for short- or medium-term use, since the body develops a tolerance to them that reduces their effectiveness and necessitates the use of progressively larger doses. Dependence on them may also occur, even in moderate dosages, and withdrawal symptoms have been observed in those who have used the drugs for only four to six weeks. In patients who have taken a benzodiazepine for many months or longer, withdrawal symptoms occur in between 15 and 40 percent of the cases and may take weeks or months to subside.

Withdrawal symptoms from benzodiazepines are of three kinds. Such severe symptoms as delirium or convulsions are rare. Frequently the withdrawal symptoms represent a renewal or increase of the anxiety itself. Many patients also experience other symptoms such as hypersensitivity to noise and light, as well as muscle twitching. As a result, many long-term users continue to take the drug not because of persistent anxiety but because the withdrawal symptoms are too unpleasant.

Because of the danger of dependence, benzodiazepines should be taken in the lowest possible dose for no more than a few weeks. For longer periods they should be taken intermittently, and only when the anxiety is severe.

Benzodiazepines act on specialized receptors in the brain that are adjacent to receptors for a neurotransmitter called gamma-aminobutyric acid (GABA), which inhibits anxiety. It is possible that the interaction of benzodiazepines with these receptors facilitates the inhibitory (anxiety-suppressing) action of GABA within the brain.

Other medications such as antidepressants or buspirone also are used in the treatment of anxiety.

Antidepressant agents. Many persons suffering from depression gain symptomatic relief from treatment with an antidepressant. There are several classes of antidepressant drugs, which vary in their mechanism of action and side effects. Successful treatment with such drugs relieves all the symptoms of depression, including disturbances of sleep and appetite, loss of sexual desire, and decreased energy, interest, and concentration. It usually takes two to three weeks for an antidepressant to improve a person's depressed mood significantly. Once a good response has been achieved, the drug should be continued for a further six months to reduce the risk of relapse. Antidepressants are also effective in treating other mental disorders such as panic disorder, agoraphobia, obsessive-compulsive disorder, and bulimia nervosa. Some of the most widely used antidepressant drugs are shown in Table 3.

It is widely theorized that depression is partly caused by reduced quantities or reduced activity of one or more neurotransmitters in the brain. Selective serotonin reuptake inhibitors (SSRIs), which include fluoxetine (Prozac) and sertraline (Zoloft), are thought to act by inhibiting the reabsorption of the neurotransmitter serotonin. As a result, there is an accumulation of serotonin in the brain, a

Side effects of benzodiazepines

Table 3: Some Commonly Used Antidepressant Drugs

class	drug	representative trade name
selective serotonin reuptake inhibitors	fluoxetine	Prozac
	paroxetine	Paxil
	sertraline	Zoloft
tricyclics	imipramine	Tofranil
	amitriptyline	Elavil
	clomipramine	Anafranil
monoamine oxidase inhibitors	phenelzine	Nardil
	tranylcypromine	Parnate
	isocarboxazid	Marplan

change that may be important in elevating mood. Because SSRIs interfere with only one neurotransmitter system, they have fewer, and less severe, side effects than other classes of antidepressants, which inhibit the action of several neurotransmitters. Common side effects include decreased sexual drive or ability, diarrhea, insomnia, headache, and nausea.

Tricyclic antidepressants (so called because of their three-ringed chemical structure) interfere with the reuptake of norepinephrine, serotonin, and dopamine. The side effects of these drugs are mostly due to their interference with the function of the autonomic nervous system and may include dryness of the mouth, blurred vision, constipation, and difficulty urinating. Weight gain can be a distressing side effect in persons taking a tricyclic for a long period of time. In elderly persons these drugs can cause delirium. Certain tricyclics interfere with conduction in heart muscle, and so they are best avoided in individuals with heart disease. Drug interactions occur with tricyclics, the most important being their interference with the action of certain drugs used in the treatment of high blood pressure.

Monoamine oxidase inhibitors (MAOIs) interfere with the action of monoamine oxidase, an enzyme involved in the breakdown of norepinephrine and serotonin. As a result, these neurotransmitters accumulate within nerve cells and presumably leak out onto receptors. The side effects of these drugs include daytime drowsiness, insomnia, and a fall in blood pressure when changing position. The MAOIs interact dangerously with various other drugs, including narcotics and some over-the-counter drugs used in treating colds. Persons taking an MAOI must avoid certain foods containing tyramine or other naturally occurring amines, which can cause a severe rise in blood pressure leading to headaches and even to stroke. Tyramine occurs in cheese, Chianti and other red wines, well-cured meats, and foods that contain monosodium glutamate (MSG).

Newer antidepressants, such as bupropion (Wellbutrin), have been introduced. These drugs are chemically unrelated to the other classes of antidepressants.

Mood-stabilizing drugs. Lithium, usually administered as its carbonate in several small doses per day, is effective in the treatment of an episode of mania. It can drastically reduce the elation, overexcitement, grandiosity, paranoia, irritability, and flights of ideas typical of people in the manic state. It has little or no effect for several days, however, and a therapeutic dose is rather close to a toxic dose. In severe episodes antipsychotic drugs may also be used. Lithium also has an antidepressant action in some patients with melancholia.

The most important use of lithium is in the maintenance treatment of patients with bipolar disorder or with recurrent depression. When given while the patient is well, lithium may prevent further mood swings, or it may reduce either their frequency or their severity. Its mode of action is unknown. Treatment begins with a small dose that is gradually increased until a specified concentration of lithium in the blood is reached. Blood tests to determine this are carried out weekly in the early stages of treatment and later every two to three months. It may take as long as a year for lithium to become fully effective.

The toxic effects of lithium, which usually occur when there are high concentrations of it in the blood, include drowsiness, coarse tremors, vomiting, diarrhea, incoordination of movement, and, with still higher blood concentrations, convulsions, coma, and death. At therapeutic blood concentrations, lithium's side effects include fine

tremors (which can be alleviated by propranolol), weight gain, passing increased amounts of urine with consequent increased thirst, and reduced thyroid function.

Carbamazepine, an anticonvulsant drug, has been shown to be effective in the treatment of mania and in the maintenance treatment of bipolar disorder. It may be combined with lithium in patients with bipolar disorder who fail to respond to either drug alone. Divalproex, another anticonvulsant, is also used in the treatment of mania.

Electroconvulsive treatment. In electroconvulsive therapy (ECT), also called shock treatment, a seizure is induced in a patient by passing a mild electric current through his brain. The mode of action of ECT is not understood. Several studies have shown that ECT is effective in treating patients with severe depression, acute mania, and some types of schizophrenia. However, the procedure remains controversial and is used only if all other methods of treatment have failed.

Prior to the administration of ECT, the patient is given an intravenous injection of an anesthetic to put him to sleep and then is administered an injection of a muscle-relaxant in order to reduce muscular contractions during the treatment. The electrical current is then applied to the brain. In bilateral ECT this is done by applying an electrode to each side of the head; in unilateral ECT both electrodes are placed over the nondominant cerebral hemisphere—*i.e.*, the right side of the head in a right-handed person. Unilateral ECT produces noticeably less confusion and memory impairment in patients, but more treatments may be needed. Patients recover consciousness rapidly after the treatment but may be confused and may experience a mild headache for an hour or two.

ECT treatments are normally given two or three times a week in the treatment of patients with depression. The number of electroconvulsive treatments required for treating depression is usually between 6 and 12. Some patients improve after the first treatment, others only after several. Once a program of ECT has been successfully completed, maintenance treatment with an antidepressant significantly decreases the patient's risk of relapse.

ECT is often considered for cases of severe depression when the patient's life is endangered because of refusal of food and fluids or because of serious risk of suicide, as well as in cases of postpartum depression, when it is desirable to reunite the mother and baby as soon as possible. ECT is often used in treating patients whose depression has not responded to adequate dosages of antidepressants.

The chief unwanted effect of ECT is impairment of memory. Some patients report memory gaps covering the period just before treatment, but others lose memories from several months before treatment. Many patients have memory difficulties for a few days or even a few weeks after completion of the treatment so that they forget appointments, phone numbers, and the like. These difficulties are transient and disappear rapidly in the vast majority of patients. Occasionally, however, patients complain of permanent memory impairment after ECT.

Psychosurgery. Psychosurgery is the destruction of groups of nerve cells or nerve fibres in the brain by surgical techniques in an attempt to relieve severe psychiatric symptoms. The removal of a brain tumour that is causing psychiatric symptoms is not an example of psychosurgery.

The classical technique of bilateral prefrontal leucotomy (lobotomy) is no longer performed because of its frequent undesirable effects on physical and mental health, in particular the development of epilepsy and the appearance of permanent undesirable changes in personality. The latter include increased apathy and passivity, lack of initiative, and a generally decreased depth and intensity of the person's emotional responses to life. In its heyday the procedure was used to treat chronically self-destructive, delusional, agitated, or violent psychotic patients. Stereotaxic surgical techniques have been developed that enable the surgeon to produce small areas of nerve cell or fibre destruction by means of metal probes inserted into accurately located parts of the brain. The nerve tissue is then destroyed by the implantation of a radioactive substance (usually yttrium) or by the application of heat or cold.

The proponents of psychosurgery claim that it is effective

in treating some patients with severe and intractable obsessive-compulsive disorder and that it may improve the behaviour of abnormally aggressive patients. Many of the therapeutic effects that were claimed for psychosurgery by its adherents are, in fact, now attainable by the use of anti-psychotic and antidepressant medications. Today psychosurgery has a very small part to play in psychiatric treatment when the prolonged use of other forms of treatment has been unsuccessful and the patient is chronically and severely distressed or tormented by his symptoms. Whereas ECT is a routine treatment in certain specified conditions, psychosurgery is, at best, a last resort.

THE PSYCHOTHERAPIES

Psychotherapy implies the treatment of individuals with mental discomfort, dysfunction, or disease by psychological means by a trained therapist who adheres to a particular theory of both symptom causation and symptom relief. The American psychotherapist Jerome Frank has classified psychotherapies into religio-magical and empirico-scientific forms. The former depend on the shared beliefs of the therapist and client in magic, spirits, or other supernatural processes or powers. This article is concerned, however, with the latter forms of psychotherapy, which have been developed by modern medicine and which are performed by a member of one of the mental health professions such as a psychiatrist or a clinical psychologist.

It is usual to contrast two main forms of psychotherapy, dynamic and behavioral. They are conceptually different; behaviour therapy concentrates on alleviating a patient's overt symptoms, which are attributed to faulty learning, while dynamic therapy concentrates on understanding the meaning of symptoms and understanding the emotional conflicts within the patient that may be causing those symptoms. In their pure forms the two approaches are very different, but in practice many therapists use elements of both.

Dynamic psychotherapies. There are many variants of dynamic psychotherapy, most of which ultimately derive from the basic precepts of psychoanalysis. The fundamental approach of most dynamic psychotherapies can be traced to three basic theoretical principles or assertions: (1) Human behaviour is prompted chiefly by emotional considerations, but insight and self-understanding are necessary to modify and control such behaviour and its underlying aims; (2) A significant proportion of human emotion is not normally accessible to one's personal awareness or introspection, being rooted in the unconscious, those portions of the mind beneath the level of consciousness; (3) Any process that makes available to a person's conscious awareness the true significance of emotional conflicts and tensions that were hitherto held in the unconscious will thereby produce heightened awareness and increased stability and emotional control. The classic dynamic psychotherapies are relatively intensive talking treatments that are aimed at providing people with insight into their conscious and unconscious mental processes, with the goal of enabling them ultimately to achieve a better self-understanding.

Dynamic psychotherapy attempts to enhance the patient's personality growth as well as to alleviate symptoms. The main therapeutic forces are activated in the relationship between patient and therapist and depend both upon the empathy, understanding, integrity, and concern demonstrated by the therapist and upon the motivation, intelligence, and capacity for achieving insight manifested by the patient. The attainment of a therapeutic alliance—i.e., a working relationship between patient and therapist that is based on mutual respect, trust, and confidence—provides the context in which the patient's problems can be worked through and resolved. Several of the most important forms are treated below.

Psychoanalytic psychotherapy. Classical psychoanalysis is the most intensive of all the psychotherapies in terms of time, cost, and effort. It is conducted with the patient lying on a couch and with the analyst seated out of sight but close enough to hear what the patient says. The treatment sessions last 50 minutes and are usually held four or five times a week for at least three years. The primary tech-

nique used in psychoanalysis and in other dynamic psychotherapies to enable unconscious material to enter the patient's consciousness is that of "free association." In free association, according to Freud, the patient "is to tell us not only what he can say intentionally and willingly, what will give him relief like a confession, but everything else as well that his self-observation yields him, everything that comes into his head, even if it is disagreeable for him to say it, even if it seems to him unimportant or actually nonsensical." Such a procedure is rendered difficult, first, because for a person to speak of his innermost (and often socially unacceptable) thoughts is a departure from years of practice in which he has selected what he has said to others. Free association is also difficult because the patient resists remembering repressed experiences or feelings that are connected with intense or conflicting emotions that have never been finally resolved or settled. Such repressed emotions or memories usually revolve around the patient's important personal relationships and his innermost feelings of self, and the release or recollection of such emotions in the course of treatment can be itself intensely disturbing.

Attentive listening and empathy on the part of the therapist allows the patient to express thoughts and feelings that in turn permit the uncovering of underlying emotional conflicts. In the course of treatment, the patient often seeks to project (attribute to something other than himself) the disturbing emotions he feels in the process of recollection and free association, and the person who is almost invariably selected for the focus of such projection is the psychoanalyst; that is, the patient is likely to blame his emotional distress on the analyst. In this way the patient comes to feel love or hatred, dependence or rebellion, and rivalry or rejection toward the analyst. These are the same attitudes the patient has felt but has never consciously acknowledged toward his parents or other people with whom he shared important relations earlier in life. The patient's projection onto the therapist of these feelings and behaviours that originated in earlier relationships is called transference. To facilitate the development of transference, the analyst endeavours to maintain a neutral stance toward the patient in order to serve as a "blank screen" onto which the patient can project his inner feelings. The analyst's handling of the transference situation is of vital importance in psychoanalysis—or, indeed, in any form of dynamic psychotherapy. It is through transference that the patient discovers the nature of his unconscious feelings and then becomes able to acknowledge them. Once this has been done, he often finds himself able to regard his inner feelings in a far more dispassionate and tolerant light and often feels liberated from their influence.

A major therapeutic tool in the course of treatment is interpretation. This technique helps the patient to become aware of previously repressed aspects of his emotional conflicts and to uncover the meaning of uncomfortable feelings evoked by transference. Interpretation, in turn, is used to determine the underlying psychological meaning of the patient's dreams, which are held to have a hidden or latent content that symbolizes and indirectly expresses aspects of the patient's emotional conflicts.

Individual dynamic psychotherapy. Psychoanalysis has had a profound influence, particularly on American psychiatry. During the 1980s and '90s, fewer patients entered psychoanalysis, and many analysts carry out short-term individual dynamic psychotherapy. This form of therapy usually requires 50 minutes a week for six to 18 months. The aim of treatment, as in psychoanalysis, is to increase the patient's insight (understanding of himself), to relieve his symptoms, and to improve his psychological functioning. Additionally, the therapist provides more support of and structure to the patient. Suitable patients include those with a wide range of psychological disorders and personal or social problems who wish to change and who are able to view their problems in psychological terms.

As in psychoanalysis, the patient learns to trust the therapist and becomes able to talk candidly and honestly about his most intimate thoughts and feelings. The treatment setting is less formal than in psychoanalysis, with the therapist and patient seated so that eye contact can be achieved if desired.

Transference

The role of unconscious mental processes

Treatment techniques include free association and the use of interpretation by the therapist to analyze transference, the patient's unconscious defense mechanisms, and his dreams. The therapist may ask the patient to clarify or enlarge on some point on which the therapist is not clear if this seems important in the development of the patient's symptoms. The therapist directs the patient's attention to important links, of which he seems unaware, between the present and the past, between his emotional responses to the therapist and to people important to him, and so on. The therapist may challenge the patient with the likely consequences of his resistant or maladaptive behaviour and stress instead the importance of confronting and trying to resolve his psychological difficulties.

Brief focal psychotherapy. This is a form of short-term dynamic therapy in which a time limit to the duration of the therapy is often agreed upon with the patient at the outset. Sessions lasting 30 to 60 minutes are held weekly for, typically, 5 to 15 weeks. At the beginning of treatment the therapist helps identify the patient's problem or problems, and these are made the focus of the treatment. The problem should be an important source of distress to the patient, and it should be modifiable within the time limit. The therapist is more active, directive, and confrontational than in long-term dynamic therapy and ensures that the patient keeps to the focus of treatment and is not diverted by subsidiary problems or concerns. Some therapists may aim to produce considerable emotional arousal in the patient during each session as a way to activate or highlight specific problems.

Group psychotherapy. Many types of psychological treatment may be provided for groups of patients with psychiatric disorders. This is true, for example, of relaxation training and anxiety-management training. There are also self-help groups, of which Alcoholics Anonymous is perhaps the best known. A considerable number of group experiences have been devised for people who are not suffering from any psychiatric disorder; encounter groups are a well-known example. This discussion, however, is concerned with long-term dynamic group therapy, in which 6 to 10 psychiatric patients meet with a trained group therapist, or sometimes two therapists, for 60 to 90 minutes a week for up to 18 months. Often the group is closed, *i.e.*, confined to the original group membership, even if one or more members drop out before the treatment ends. In an open group patients who have stopped attending, whether by default or because of the relief of symptoms, are replaced by new members.

The types of mental disorders considered suitable for group therapy are much the same as those suitable for individual therapy. Again the patients must want to change and must be psychologically minded. In addition, it is important that they not consider group therapy as a poor second to individual therapy.

There are many varieties of dynamic group therapy, and they differ in their theoretical background and technique. The influential model of the American psychiatrist Irvin D. Yalom provides a good example of such therapies. The therapist continually encourages the patients to direct their attention to the personal interactions occurring within the group rather than to what happened in the past to individual members or what is currently happening outside the group, although both of these areas may be considered when they are relevant. The therapist regularly draws attention to what is happening between members of the group as they learn more about themselves and test out different ways of behaving with one another. The goal in group therapy is to create a climate in which the participants can shed their inhibitions. When the members come to trust one another, they are able to provide feedback and to respond to other group members in ways that might not be possible in ordinary social interactions because of the constraints of social conventions.

Several factors appear to be important in group therapy. The most important is group cohesion, which gives the patient a feeling of belonging, identification, and security and enables him to be frank and open and to take risks without the danger of rejection. Universality refers to the patient's realization that he is not unique, that all the other

group members have problems, some of them similar to his. Optimism about what can be achieved in the group, fostered by the perception of change in others, combats demoralization. Guidance, the giving of advice and explanation, is important in the early meetings of the group and is largely a function of the therapist. What has been called vicarious learning later becomes more important; the patient observes how other group members evolve solutions to common problems and emulates desirable qualities he sees in fellow members. Catharsis, or the release of highly charged emotion, occurs within the group setting and is helpful, provided that the patient can come to understand it and appreciate its significance. Another factor that is helpful in improving self-esteem is altruism, the opportunity to give assistance to another group member.

Family therapy. Family therapists view the family as the "patient" or "client" and as more than the sum of its members. The family as a focus for treatment usually comprises the members who live under the same roof, sometimes supplemented by relatives who live elsewhere or by non-relatives who share the family home. Therapy with couples—marital therapy—may be considered as a special type of family therapy. Family therapy may be appropriate when the person referred for treatment has symptoms clearly related to such disturbances in family function as marital discord, distorted family roles, and parent-child conflict or when the family as a unit asks for help. It is not appropriate when the patient has a severe disorder needing specific treatment in its own right.

The many theoretical approaches include psychoanalytic, systems-theory, and behavioral models. The analyst is concerned with the family's past as the cause of the present and pays attention to psychodynamic aspects of the individual members as well as of the family unit. The analyst also makes numerous interpretations and aims at increasing the insight of the members. The systems therapist, on the other hand, is interested in the present rather than the past and is often concerned not with promoting insight but rather with changing the family system, perhaps by altering the implicit and fixed rules under which it functions so that it can do so more effectively. The behaviour therapist is concerned with behaviour patterns and with pinpointing the types of reinforcement that maintain behaviour that other family members regard as undesirable. Members specify the changes in behaviour that they wish to see in one another, and strategies are devised to reinforce the desired behaviours. This approach has been shown to be effective in work with couples, when one partner promises some particular change, provided that the other reciprocates.

Treatment sessions in family therapy are rarely held more often than once a week and often take place only once every three or four weeks. Termination commonly occurs when the therapist considers that treatment has succeeded—or failed irretrievably—or the family decides to withdraw from treatment. There seems no doubt that family therapy can produce marked change within a family.

Behavioral psychotherapy. This is an approach to the treatment of mental disorders that uses a variety of methods based upon principles derived from experimental psychology, mainly that of learning theory. In the words of Joseph Wolpe, "Behavior therapy, or conditioning therapy, is the use of experimentally established principles of learning for the purpose of changing unadaptive behavior. Unadaptive habits are weakened and eliminated; adaptive habits are initiated and strengthened" (*The Practice of Behavior Therapy*, 1973).

In the treatment of phobias, behaviour therapy seeks to modify and eliminate the avoidance response that the patient manifests when confronted with a phobic object or situation. This is crucial because, since the patient's avoidance of the anxiety-producing situation does indeed reduce his anxiety, his conditioned association of the phobic situation with the experience of anxiety remains unchallenged and frequently persists. Behaviour therapy interrupts this self-reinforcing pattern of avoidance behaviour by presenting the feared situation to the patient in a controlled manner such that it eventually ceases to produce anxiety; the strong associative link that has been built up within

Treatment
of phobias

him between the feared situation, the experience of anxiety, and his subsequent avoidance behaviour will thus be broken down and replaced by a less maladaptive set of responses.

The behavioral therapist is concerned with the forces and mechanisms that perpetuate the patient's present symptoms or abnormal behaviours, not with experiences in the past that may have caused them or with any postulated intrapsychic conflict. Behavioral therapy focuses on observable phenomena—*i.e.*, what is done and what is said, rather than on what must be inferred (unconscious motives and processes and symbolic meanings).

The behavioral therapist carries out a detailed analysis of the patient's behaviour problems, paying particular attention to the circumstances in which they occur, to the patient's attempts to cope with his symptoms, and to his wish for change. The goals of treatment are precisely defined and usually do not include aims such as personal growth or personality change. The relationship between patient and therapist is sometimes said to be unimportant in behaviour therapy, and it is quite true that a patient may successfully follow a behavioral therapeutic program from a book or a personal computer. Nevertheless, a patient is more likely to complete an arduous program if he trusts and respects the therapist.

Behaviour therapy has become the preferred treatment for phobic states and for some obsessive-compulsive disorders, and it is effective in many cases of sexual dysfunction and deviation. It also has an important role in the rehabilitation of patients with chronic disabling disorders.

The essence of the treatment of phobias is the controlled exposure of the patient to the objects or situations that he fears. Behaviour therapy tries to eliminate the phobia by teaching the patient how to face those situations that clearly trigger his discomfort so that he can learn to tolerate them. The exposure of the patient to the feared situation can be gradual (sometimes called desensitization) or rapid (sometimes known as flooding). Contrary to popular belief, the anxiety that is produced during exposure is not usually harmful. Even if severe panic does strike the sufferer, it will gradually evaporate and will be less likely to return in the future. Effective exposure treatments developed only as therapists learned to endure the phobic anxiety exhibited by their patients and when the therapists grew secure in the knowledge that such anxiety is much more likely to lead to the patients' improvement than to be harmful. The important point in this therapy is to persevere until the phobic anxiety starts to lessen and to be prepared to go on until it does. In general, the more rapidly and directly the worst fears are embraced by the patient, the more quickly the phobic terror fades to a tolerable mild tension.

In the technique of desensitization, the patient is first taught how to practice muscular relaxation. He then reviews the situations of which he is afraid and lists them in order of increasing fearfulness, called a "hierarchy." Finally, the patient faces the various fear-producing situations in ascending order by means of vividly imagining them, countering any anxiety by using relaxation techniques. This treatment is prolonged, and its use is restricted to feared situations that patients cannot regularly confront in real life, such as fear of lightning.

One of the most common phobic disorders treated by exposure techniques is agoraphobia (fear of open or public places). The patient is encouraged to practice exposure daily, staying in a phobic situation for at least an hour, so that anxiety has time to reach a peak and then subside. The patient must be determined to get the better of the fears and not to run away from them. People with agoraphobia cannot simply enter a dreaded crowded store, feel the familiar surge of panic, and rush out again. They must devote a full afternoon to a shopping trip. When panic strikes, the patient can sit in a corner of the store and ride out the terror. When feeling better, the patient can continue shopping. Persistence and patience are essential to conquering phobias in this way.

There is considerable evidence that exposure techniques work in most cases. Even phobias present for as long as 20 years can be overcome in a treatment program requir-

ing no more than 3 to 15 hours of therapist time per patient. There is also considerable evidence that many people with phobias can treat themselves perfectly adequately by self-exposure without a therapist, using carefully devised self-help manuals.

Some patients with obsessive-compulsive disorders can also be helped by behaviour therapy. Several different techniques may be needed. For instance, a patient with an obsessional fear of contamination is treated by exposure, being taught to "soil" his hands with dirt and then to avoid washing them for longer and longer periods. Anxiety-management training enables the patient to withstand the anxiety he feels during the period of exposure.

This and other techniques have been shown to be effective in the treatment of compulsive rituals, with improvement occurring in more than two-thirds of patients. There is also a reduction in the frequency and intensity of obsessional thoughts that accompany the rituals. The treatment of obsessional thoughts that occur alone is much less satisfactory, however.

Other therapies. Many other types of psychotherapy were developed in the second half of the 20th century, each with its own emphasis on symptom causation and its own particular approach to treatment. Many of these therapies use classical dynamic and behavioral models in modified forms, and they may also stress the understanding and modification of cognition and the ways in which people "process" their experiences, moods, and emotions. Among these relatively recent psychotherapies are client-centred psychotherapy, developed by the American psychologist Carl R. Rogers; transactional analysis, originated by the American psychiatrist Eric Berne; the interpersonal therapies developed by the American psychiatrists Adolf Meyer and Harry Stack Sullivan; cognitive therapy, developed by the American psychiatrist Aaron T. Beck; rational-emotive psychotherapy, developed by the American psychologist Albert Ellis; and Gestalt therapy, which stems from the work of the German psychiatrist Frederick S. ("Fritz") Perls.

Another class of therapies consists of those used to care for psychotic patients, both those in hospitals and those who live in the community. Supportive psychotherapy consists of the long-term treatment of patients with chronic schizophrenia or other mental disorders. Such a therapy uses reassurance, guidance, and encouragement to help the patient cope with his disabilities and live as satisfactory a life as possible. Rehabilitation programs for chronic or episodically psychotic patients include medication maintenance; training in social skills that they may have lost while ill; occupational training to improve the patient's skills in cooking, shopping, and other domestic tasks; and industrial therapy, which usually offers the patient gainful employment under conditions of minimal stress. Family therapy is sometimes used to help relatives learn to cope with a family member who has schizophrenia and has returned home from the hospital.

Community care for individuals with schizophrenia or other psychoses who have been discharged from a hospital must provide them with medication maintenance and a minimum of psychiatric monitoring, appropriate housing facilities, some type of employment, and training in such skills as using public transport, preparing their own food, and looking after their finances. Each patient should have a case manager, a professional social worker who maintains contact and secures from governmental or social agencies the assistance that the patient needs. When provisions such as these are not made, many formerly hospitalized patients stop taking their medicine and, in effect, drop out of the mental health care system, becoming unemployed and even homeless. This phenomenon became particularly evident in the United States and to a lesser extent in western Europe when massive numbers of mental patients were released from hospitals during the 1950s and '60s after the effectiveness of antipsychotic medications had been verified. These releases were also motivated by the concerns of civil libertarians over the abuse of patients' rights in keeping them committed to mental hospitals. However, the support network of community-based mental health clinics that would have been necessary to cope

Unforeseen effects of deinstitutionalization

with the released patients was either inadequately established or nonexistent. The result was that many psychotic patients received inadequate outpatient care and supervision or encountered severe difficulties in obtaining housing or employment, becoming homeless wanderers in large urban areas.

(J.L.Gi./S.C.Y.)

BIBLIOGRAPHY

General works: The following works provide descriptions of the syndromes, causes, epidemiology, and methods of treatment of mental disorders: HAROLD I. KAPLAN and BENJAMIN J. SADOCK, *Kaplan and Sadock's Synopsis of Psychiatry: Behavioral Sciences, Clinical Psychiatry*, 8th ed. (1998); ROBERT E. HALES and STUART C. YUDOFKY (eds.), *Essentials of Clinical Psychiatry* (1999), based on the following: ROBERT E. HALES, STUART C. YUDOFKY, and JOHN A. TALBOTT (eds.), *The American Psychiatric Press Textbook of Psychiatry*, 3rd ed. (1999); STUART C. YUDOFKY and ROBERT E. HALES (eds.), *The American Psychiatric Press Textbook of Neuropsychiatry*, 3rd ed. (1997); and HAROLD I. KAPLAN and BENJAMIN J. SADOCK (eds.), *Comprehensive Textbook of Psychiatry/VII*, 6th ed., 2 vol. (1995).

Classification and epidemiology: The two classificatory systems mentioned in the text are detailed in AMERICAN PSYCHIATRIC ASSOCIATION, *Diagnostic and Statistical Manual of Mental Disorders: DSM-IV*, 4th ed. (1994, reissued 1998); and WORLD HEALTH ORGANIZATION, *International Statistical Classification of Diseases and Related Health Problems*, 10th revision, 3 vol.

Causation: Different aspects and theories of causation are considered in GLEN O. GABBARD, *Psychodynamic Psychiatry in Clinical Practice*, 3rd ed. (2000); ROBERT B. WHITE and ROBERT M.

GILLILAND, *Elements of Psychopathology: The Mechanisms of Defense* (1975); FREDERICK K. GOODWIN and KAY REDFIELD JAMISON, *Manic-Depressive Illness* (1990); and NANCY C. ANDREASEN (ed.), *Schizophrenia: From Mind to Molecule* (1994).

Treatment: A comprehensive discussion of current treatment methods for various mental disorders is offered in ROBERT E. HALES, STUART C. YUDOFKY, and JOHN A. TALBOTT (eds.), *The American Psychiatric Association Textbook of Psychiatry*, 3rd ed. (1999). A work concerned with the theoretical concepts underlying psychotherapy is CHARLES BRENNER, *An Elementary Textbook of Psychoanalysis*, rev. ed. (1973, reissued 1990). A good account of the main forms of psychological treatment is given in SIDNEY BLOCH (ed.), *An Introduction to the Psychotherapies*, 3rd ed. (1996). Also of interest are PETER E. SIFNEOS, *Short-Term Dynamic Psychotherapy: Evaluation and Technique*, 2nd ed. (1987); IRVIN D. YALOM, *The Theory and Practice of Group Psychotherapy*, 4th ed. (1995); and JOHN C. MASTERS et al., *Behavior Therapy: Techniques and Empirical Findings*, 3rd ed. (1987).

Pharmacological and physical methods of treatment are dealt with in ALAN F. SCHATZBERG and CHARLES B. NEMEROFF (eds.), *The American Psychiatric Press Textbook of Psychopharmacology*, 2nd ed. (1998); ROSS J. BALDESSARINI, *Chemotherapy in Psychiatry: Principles and Practice*, rev. and enlarged ed. (1985); STEVEN E. HYMAN and ERIC J. NESTLER, *The Molecular Foundations of Psychiatry* (1993); and STUART C. YUDOFKY, ROBERT E. HALES, and TOM FERGUSON, *What You Need to Know About Psychiatric Drugs* (1991). ELLIOT S. VALENSTEIN, *Great and Desperate Cures: The Rise and Decline of Psychosurgery and Other Radical Treatments for Mental Illness* (1986), presents a history of the methods and personalities involved.

(J.L.Gi./A.C.P.S./L.B.An.)

The History of Ancient Mesopotamia

Mesopotamia is a historical region in southwest Asia where the world's earliest civilization developed. The name comes from a Greek word meaning "between rivers," referring to the land between the Tigris and Euphrates rivers, but the region can be broadly defined to include the area that is now eastern Syria, southeastern Turkey, and most of Iraq. This region was the centre of a culture whose influence extended throughout the Middle East and as far as the Indus valley, Egypt, and the Mediterranean. This article covers the history of Mesopotamia from the prehistoric period up to the

Arab conquest in the 7th century AD. For the history of the region in the succeeding periods, see the *Macropædia* article IRAQ. (Ed.)

For a discussion of the religions of ancient Mesopotamia, see the article MIDDLE EASTERN RELIGIONS: *Mesopotamian religions*. For a discussion of Mesopotamian visual arts, see MIDDLE EASTERN ARTS AND ARCHITECTURE, ANCIENT: *Mesopotamian arts and architecture*. For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, section 911, and the *Index*.

This article is divided into the following sections:

Mesopotamia to the end of the Old Babylonian period	860	Assyria and Babylonia from c. 1000 to c. 750 BC	879
The origins of Mesopotamian history	860	Assyria and Babylonia until Ashurnasirpal II	
The background		Shalmaneser III and Shamshi-Adad V of Assyria	
The emergence of Mesopotamian civilization		Adad-nirari III and his successors	
Sumerian civilization	866	The neo-Assyrian empire (746–609)	880
The Sumerians to the end of the Early Dynastic period		Tiglath-pileser III and Shalmaneser V	
Sumer and Akkad from 2350 to 2000 BC		Sargon II (721–705) and Marduk-apal-iddina	
The 3rd dynasty of Ur		of Babylonia	
The Old Babylonian period	872	Sennacherib	
Isin and Larsa		Esarhaddon	
Early history of Assyria		Ashurbanipal (668–627) and Shamash-shum-ukin	
The Old Babylonian empire		(668–648)	
The Hurrians		Decline of the Assyrian empire	
Mesopotamia to the end of the Achaemenian period	876	The Neo-Babylonian empire	885
The Kassites, the Mitanni, and the rise of Assyria	876	Nebuchadnezzar II	
The Kassites in Babylonia		The last kings of Babylonia	
The Hurrian and Mitanni kingdoms		Mesopotamia under the Persians	886
The rise of Assyria		Mesopotamia from c. 320 BC to c. AD 620	886
Assyria and Babylonia at the end of the 2nd millennium	878	The Seleucid period	887
Babylonia under the 2nd dynasty of Isin		The Parthian period	888
Assyria between 1200 and 1000 BC		The Sāsānian period	890
		Bibliography	892

Mesopotamia to the end of the Old Babylonian period

THE ORIGINS OF MESOPOTAMIAN HISTORY

The background. In the narrow sense, Mesopotamia is the area between the Euphrates and Tigris rivers, north or northwest of the bottleneck at Baghdad, in modern Iraq; it is Al-Jazīrah ("The Island") of the Arabs. South of this lies Babylonia, named after the city of Babylon. However, in the broader sense, the name Mesopotamia has come to be used for the area bounded on the northeast by the Zagros Mountains and on the southwest by the edge of the Arabian Plateau and stretching from the Persian Gulf in the southeast to the spurs of the Anti-Taurus Mountains in the northwest. Only from the latitude of Baghdad do the Euphrates and Tigris truly become twin rivers, the *rāfidān* of the Arabs, which have constantly changed their courses over the millennia. The low-lying plain of the Kārūn River in Persia has always been closely related to Mesopotamia, but it is not considered part of Mesopotamia as it forms its own river system.

Mesopotamia, south of Ar-Ramādī (about 70 miles, or 110 kilometres, west of Baghdad) on the Euphrates and the bend of the Tigris below Sāmarrā' (about 70 miles north-northwest of Baghdad), is flat alluvial land. Between Baghdad and the mouth of the Shaṭṭ al-'Arab (the confluence of the Tigris and Euphrates, where it empties into the Persian Gulf) there is a difference in height of only about 100 feet (30 metres). As a result of the slow flow of the water, there are heavy deposits of silt, and the riverbeds are raised. Consequently, the rivers often overflow their banks (and may even change their course) when they are not protected by high dikes. In recent times they have been regulated above Baghdad by the use of escape channels with overflow reservoirs. The extreme south is

a region of extensive marshes and reed swamps, *hawrs*, which, probably since early times, have served as an area of refuge for oppressed and displaced peoples. The supply of water is not regular; as a result of the high average temperatures and a very low annual rainfall, the ground of the plain of latitude 35° N is hard and dry and unsuitable for plant cultivation for at least eight months in the year. Consequently, agriculture without risk of crop failure, which seems to have begun in the higher rainfall zones and in the hilly borders of Mesopotamia in the 10th millennium BC, began in Mesopotamia itself, the real heart of the civilization, only after artificial irrigation had been invented, bringing water to large stretches of territory through a widely branching network of canals. Since the ground is extremely fertile and, with irrigation and the necessary drainage, will produce in abundance, southern Mesopotamia became a land of plenty that could support a considerable population. The cultural superiority of north Mesopotamia, which may have lasted until about 4000 BC, was finally overtaken by the south when the people there had responded to the challenge of their situation.

The present climatic conditions are fairly similar to those of 8,000 years ago. An English survey of ruined settlements in the area 30 miles around ancient Hatra (180 miles northwest of Baghdad) has shown that the southern limits of the zone in which agriculture is possible without artificial irrigation has remained unchanged since the first settlement of Al-Jazīrah.

The availability of raw materials is a historical factor of great importance, as is the dependence on those materials that had to be imported. In Mesopotamia, agricultural products and those from stock breeding, fisheries, date palm cultivation, and reed industries—in short, grain, vegetables, meat, leather, wool, horn, fish, dates, and reed and plant-fibre products—were available in plenty and

could easily be produced in excess of home requirements to be exported. There are bitumen springs at Hit (90 miles northwest of Baghdad) on the Euphrates (the Is of Herodotus). On the other hand, wood, stone, and metal were rare or even entirely absent. The date palm—virtually the national tree of Iraq—yields a wood suitable only for rough beams and not for finer work. Stone is mostly lacking in southern Mesopotamia, although limestone is quarried in the desert about 35 miles to the west and “Mosul marble” is found not far from the Tigris in its middle reaches. Metal can only be obtained in the mountains, and the same is true of precious and semiprecious stones. Consequently, southern Mesopotamia in particular was destined to be a land of trade from the start. Only rarely could “empires” extending over a wider area guarantee themselves imports by plundering or by subjecting neighbouring regions.

The raw material that epitomizes Mesopotamian civilization is clay: in the almost exclusively mud-brick architecture and in the number and variety of clay figurines and pottery artifacts, Mesopotamia bears the stamp of clay as does no other civilization, and nowhere in the world but in Mesopotamia and the regions over which its influence was diffused was clay used as the vehicle for writing. Such phrases as cuneiform civilization, cuneiform literature, and cuneiform law can apply only where people had had the idea of using soft clay not only for bricks and jars and for the jar stoppers on which a seal could be impressed as a mark of ownership but also as the vehicle for impressed signs to which established meanings were assigned—an intellectual achievement that amounted to nothing less than the invention of writing.

The character and achievements of ancient Mesopotamia. Questions as to what ancient Mesopotamian civilization did and did not accomplish, how it influenced its neighbours and successors, and what its legacy has transmitted are posed from the standpoint of 20th-century civilization and are in part coloured by ethical overtones, so that the answers can only be relative. Modern scholars assume the ability to assess the sum total of an “ancient Mesopotamian civilization”; but, since the publication of an article by the Assyriologist Benno Landsberger on “Die Eigenbegrifflichkeit der babylonischen Welt” (1926; “The Distinctive Conceptuality of the Babylonian World”), it has become almost a commonplace to call attention to the necessity of viewing ancient Mesopotamia and its civilization as an independent entity.

Ancient Mesopotamia had many languages and cultures; its history is broken up into many periods and eras; it had no real geographic unity, and above all no permanent capital city, so that by its very variety it stands out from other civilizations with greater uniformity, particularly that of Egypt. The script and the pantheon constitute the unifying factors, but in these also Mesopotamia shows its predilection for multiplicity and variety. Written documents were turned out in quantities, and there are often many copies of a single text. The pantheon consisted of more than 1,000 deities, even though many divine names may apply to different manifestations of a single god. During 3,000 years of Mesopotamian civilization, each century gave birth to the next. Thus classical Sumerian civilization influenced that of the Akkadians, and the Ur III empire, which itself represented a Sumero-Akkadian synthesis, exercised its influence on the first quarter of the 2nd millennium BC. With the Hittites, large areas of Anatolia were infused with the culture of Mesopotamia from 1700 BC onward. Contacts, via Mari, with Ebla in Syria, some 30 miles south of Aleppo, go back to the 24th century BC, so that links between Syrian and Palestinian scribal schools and Babylonian civilization during the Amarna period (14th century BC) may have had much older predecessors. At any rate, the similarity of certain themes in cuneiform literature and the Old Testament, such as the story of the Flood or the motif of the righteous sufferer, is due to such early contacts and not to direct borrowing.

The world of mathematics and astronomy owes much to the Babylonians—for instance, the sexagesimal system for the calculation of time and angles, which is still practical because of the multiple divisibility of the number 60; the

Greek day of 12 “double-hours”; and the zodiac and its signs. In many cases, however, the origins and routes of borrowings are obscure, as in the problem of the survival of ancient Mesopotamian legal theory.

The achievement of the civilization itself may be expressed in terms of its best points—moral, aesthetic, scientific, and, not least, literary. Legal theory flourished and was sophisticated early on, being expressed in several collections of legal decisions, the so-called codes, of which the best-known is the Code of Hammurabi. Throughout these codes recurs the concern of the ruler for the weak, the widow, and the orphan—even if, sometimes, the phrases were regrettably only literary clichés. The aesthetics of art are too much governed by subjective values to be assessed in absolute terms, yet certain peaks stand out above the rest, notably the art of Uruk IV, the seal engraving of the Akkad period, and the relief sculpture of Ashurbanipal. Nonetheless, there is nothing in Mesopotamia to match the sophistication of Egyptian art. Science the Mesopotamians had, of a kind, though not in the sense of Greek science. From its beginnings in Sumer before the middle of the 3rd millennium BC, Mesopotamian science was characterized by endless, meticulous enumeration and ordering into columns and series, with the ultimate ideal of including all things in the world but without the wish or ability to synthesize and reduce the material to a system. Not a single general scientific law has been found, and only rarely has the use of analogy been found. Nevertheless, it remains a highly commendable achievement that Pythagoras’ law (that the sum of the squares on the two shorter sides of a right-angled triangle equals the square on the longest side), even though it was never formulated, was being applied as early as the 18th century BC. Technical accomplishments were perfected in the building of the ziggurats (temple towers resembling pyramids), with their huge bulk, and in irrigation, both in practical execution and in theoretical calculations. At the beginning of the 3rd millennium BC, an artificial stone often regarded as a forerunner of concrete was in use at Uruk (160 miles south-southeast of modern Baghdad), but the secret of its manufacture apparently was lost in subsequent years.

Writing pervaded all aspects of life and gave rise to a highly developed bureaucracy—one of the most tenacious legacies of the ancient Middle East. Remarkable organizing ability was required to administer huge estates, in which, under the 3rd dynasty of Ur, for example, it was not unusual to prepare accounts for thousands of cattle or tens of thousands of bundles of reeds. Similar figures are attested at Ebla, three centuries earlier.

Above all, the literature of Mesopotamia is one of its finest cultural achievements. Though there are many modern anthologies and chrestomathies (compilations of useful learning), with translations and paraphrases of Mesopotamian literature, as well as attempts to write its history, it cannot truly be said that “cuneiform literature” has been resurrected to the extent that it deserves. There are partly material reasons for this: many clay tablets survive only in a fragmentary condition, and duplicates that would restore the texts have not yet been discovered, so that there are still large gaps. A further reason is the inadequate knowledge of the languages: insufficient acquaintance with the vocabulary and, in Sumerian, major difficulties with the grammar. Consequently, another generation of Assyriologists will pass before the great myths, epics, lamentations, hymns, “law codes,” wisdom literature, and pedagogical treatises can be presented to the reader in such a way that he can fully appreciate the high level of literary creativity of those times.

The classical and medieval views of Mesopotamia; its re-discovery in modern times. Before the first excavations in Mesopotamia, about 1840, nearly 2,000 years had passed during which knowledge of the ancient Middle East was derived from three sources only: the Bible, Greek and Roman authors, and the excerpts from the writings of Berosus, a Babylonian who wrote in Greek. In 1800 very little more was known than in AD 800, although these sources had served to stir the imagination of poets and artists, down to *Sardanapalus* (1821) by the 19th-century English poet Lord Byron.

Scarcity of wood, stone, and metal

Influence on surrounding areas

Mesopotamian literature

Greek
accounts of
Mesopotamia

Apart from the building of the Tower of Babel, the Old Testament mentions Mesopotamia only in those historical contexts in which the kings of Assyria and Babylonia affected the course of events in Israel and Judah: in particular Tiglath-pileser III, Shalmaneser V, and Sennacherib, with their policy of deportation, and the Babylonian Exile introduced by Nebuchadrezzar II. Of the Greeks, Herodotus of Halicarnassus (5th century BC, a contemporary of Xerxes I and Artaxerxes I) was the first to report on "Babylon and the rest of Assyria"; at that date the Assyrian empire had been overthrown for more than 100 years. The Athenian Xenophon took part in an expedition (during 401–399 BC) of Greek mercenaries who crossed Anatolia, made their way down the Euphrates as far as the vicinity of Baghdad, and returned up the Tigris after the famous Battle of Cunaxa. In his *Cyropaedia* Xenophon describes the final struggle between Cyrus II and the Neo-Babylonian empire. Later, the Greeks adopted all kinds of fabulous tales about King Ninus, Queen Semiramis, and King Sardanapalus. These stories are described mainly in the historical work of Diodorus Siculus (1st century BC), who based them on the reports of a Greek physician, Ctesias (405–359 BC). Herodotus saw Babylon with his own eyes, and Xenophon gave an account of travels and battles. All later historians, however, wrote at second or third hand, with one exception, Berosus (b. c. 340 BC), who emigrated at an advanced age to the Aegean island of Cos, where he is said to have composed the three books of the *Babylōniaka*. Unfortunately, only extracts from them survive, prepared by one Alexander Polyhistor (1st century BC), who, in his turn, served as a source for the Church Father Eusebius (d. AD 342). Berosus derided the "Greek historians" who had so distorted the history of his country. He knew, for example, that it was not Semiramis who founded the city of Babylon, but he was himself the prisoner of his own environment and cannot have known more about the history of his land than was known in Babylonia itself in the 4th century BC.

Berosus' first book dealt with the beginnings of the world and with a myth of a composite being, Oannes, half fish, half man, who came ashore in Babylonia at a time when men still lived like the wild beasts. Oannes taught them the essentials of civilization: writing, the arts, law, agriculture, surveying, and architecture. The name Oannes must have been derived from the cuneiform U'anna (Sumerian) or Umanna (Akkadian), a second name of the mythical figure Adapa, the bringer of civilization. The second book of Berosus contained the Babylonian king list from the beginning to King Nabonassar (Nabu-našir, 747–734 BC), a contemporary of Tiglath-pileser III. Berosus' tradition, beginning with a list of primeval kings before the Flood, is a reliable one; it agrees with the tradition of the Sumerian king list, and even individual names can be traced back exactly to their Sumerian originals. Even the immensely long reigns of the primeval kings, which lasted as long as "18 sars" (= $18 \times 3,600 = 64,800$) of years, are found in Berosus. Furthermore, he was acquainted with the story of the Flood, with Cronus as its instigator and Xisuthros (or Ziusudra) as its hero, and with the building of an ark. The third book is presumed to have dealt with the history of Babylonia from Nabonassar to the time of Berosus himself.

Diodorus made the mistake of locating Nineveh on the Euphrates, and Xenophon gave an account of two cities, Larissa (probably modern Nimrūd [ancient Kalakh], 20 miles southeast of modern Mosul) and Mespila (ancient Nineveh, just north of Mosul). The name Mespila probably was nothing more than the word of the local Aramaeans for ruins; there can be no clearer instance of the rift that had opened between the ancient Middle East and the classical West. In sharp contrast, the East had a tradition that the ruins opposite Mosul (in north Iraq) concealed ancient Nineveh. When a Spanish rabbi from Navarre, Benjamin of Tudela, was traveling in the Middle East between 1160 and 1173, Jews and Muslims alike knew the position of the grave of the prophet Jonah. The credit for the rediscovery of the ruins of Babylon goes to an Italian, Pietro della Valle, who correctly identified the vast ruins north of modern Al-Hillah, Iraq (60 miles south

The work
of Pietro
della Valle

of Baghdad); he must have seen there the large rectangular tower that represented the ancient ziggurat. Previously, other travelers had sought the Tower of Babel in two other monumental ruins: Birs Nimrūd, the massive brick structure of the ziggurat of ancient Borsippa (modern Birs, near Al-Hillah), vitrified by lightning, and the ziggurat of the Kassite capital, Dur Kurigalzu, at Burj 'Aqarqūf, 22 miles west of Baghdad. Pietro della Valle brought back to Europe the first specimens of cuneiform writing, stamped brick, of which highly impressionistic reproductions were made. Thereafter, European travelers visited Mesopotamia with increasing frequency, among them Carsten Niebuhr (an 18th-century German traveler), Claudius James Rich (a 19th-century Orientalist and traveler), and Ker Porter (a 19th-century traveler).

In modern times a third Middle Eastern ruin drew visitors from Europe—Persepolis, in the land of Persia east of Susiana, near modern Shirāz, Iran. In 1602, reports had filtered back to Europe of inscriptions that were not in Hebrew, Arabic, Aramaic, Georgian, or Greek. In 1700 an Englishman, Thomas Hyde, coined the term "cuneiform" for these inscriptions, and by the middle of the 18th century it was known that the Persepolis inscriptions were related to those of Babylon. Niebuhr distinguished three separate alphabets (Babylonian, Elamite, and Old Persian cuneiform). The first promising attempt at decipherment was made by the German philologist Georg Friedrich Grotefend in 1802, by use of the kings' names in the Old Persian versions of the trilingual inscriptions, although his later efforts led him up a blind alley. Thereafter, the efforts to decipher cuneiform gradually developed in the second half of the 19th century into a discipline of ancient Oriental philology, which was based on results established through the pioneering work of Emile Burnouf, Edward Hincks, Sir Henry Rawlinson, and many others.

Today this subject is still known as Assyriology, because at the end of the 19th century the great majority of cuneiform texts came from the Assyrian city of Nineveh, in particular from the library of King Ashurbanipal in the mound of Kuyunjik at Nineveh.

Modern archaeological excavations. More than 150 years separate the first excavations in Mesopotamia—adventurous expeditions involving great personal risks, far from the protection of helpful authorities—from those of the present day with their specialist staffs, modern technical equipment, and objectives wider than the mere search for valuable antiquities. The progress of six generations of excavators has led to a situation in which less is recovered more accurately; in other words, the finds are observed, measured, and photographed as precisely as possible. At first digging was unsystematic, with the consequence that, although huge quantities of clay tablets and large and small antiquities were brought to light, the locations of the finds were rarely described with any accuracy. Not until the beginning of the 20th century did excavators learn to isolate the individual bricks in the walls that had previously been erroneously thought to be nothing more than packed clay; the result was that various characteristic brick types could be distinguished and successive architectural levels established. Increased care in excavation does, of course, carry with it the risk that the pace of discovery will slow down. Moreover, the eyes of the local inhabitants are now sharpened and their appetite for finds is whetted, so that clandestine diggers have established themselves as the unwelcome colleagues of the archaeologists.

A result of the technique of building with mud brick (mass production of baked bricks was impossible because of the shortage of fuel) was that the buildings were highly vulnerable to the weather and needed constant renewal; layers of settlement rapidly built up, creating a tell (Arabic: *tall*), a mound of occupation debris that is the characteristic ruin form of Mesopotamia. The word itself appears among the most original vocabulary of the Semitic languages and is attested as early as the end of the 3rd millennium BC. Excavation is made more difficult by this mound formation, since both horizontal and vertical axes have to be taken into account. Moreover, the depth of each level is not necessarily constant, and foundation trenches may be dug down into earlier levels. A further

The tell

problem is that finds may have been removed from their original context in antiquity. Short-lived settlements that did not develop into mounds mostly escape observation, but aerial photography can now pick out ground discolorations that betray the existence of settlements. Districts with a high water level today, such as the reed marshes (*hawrs*), or ruins that are covered by modern settlements, such as Irbil (ancient Arbela), some 200 miles north of Baghdad, or sites that are surmounted by shrines and tombs of holy men are closed to archaeological research.

Excavations in Mesopotamia have mostly been national undertakings (France, England, the United States, Germany, Iraq, Denmark, Belgium, Italy, Japan, and the former Soviet Union), but joint expeditions like the one sent to Ur (190 miles south-southeast of Baghdad) in the 1920s have become more frequent since the 1970s. The history of archaeological research in Mesopotamia falls into four categories, represented by phases of differing lengths: the first, and by far the longest, begins with the French expedition to Nineveh (1842) and Khorsabad (the ancient Dur Sharrukin, 20 miles northeast of modern Mosul; 1843–55) and that of the English to Nineveh (1846–55) and Nimrud (ancient Kalakh, biblical Calah; 1845, with interruptions until 1880). This marked the beginning of the “classic” excavations in the important ancient capitals, where spectacular finds might be anticipated. The principal gains were the Assyrian bull colossi and wall reliefs and the library of Ashurbanipal from Nineveh, although the ground plans of temples and palaces were quite as valuable. While these undertakings had restored the remains of the Neo-Assyrian empire of the 1st millennium BC, from 1877 onward new French initiatives in Telloh (Arabic: Tall Lōh), 155 miles southeast of Baghdad, reached almost 2,000 years further back into the past. There they rediscovered a people whose language had already been encountered in bilingual texts from Nineveh—the Sumerians. Telloh (ancient Girsu) yielded not only inscribed material that, quite apart from its historical interest, was critical for the establishment of the chronology of the second half of the 3rd millennium BC but also many artistic masterpieces. Thereafter excavations in important cities spread to form a network including Susa, 150 miles west of Esfahan in Iran (France; 1884 onward); Nippur, 90 miles southeast of Baghdad (the United States; 1889 onward); Babylon, 55 miles south of Baghdad (Germany; 1899–1917 and again from 1957 onward); Ashur, modern Ash-Sharqāt, 55 miles south of Mosul (Germany; 1903–14); Uruk (Germany; 1912–13 and from 1928 onward); and Ur (England and the United States; 1918–34). Mention also should be made of the German excavations at Boğazköy in central Turkey, the ancient Hattusa, capital of the Hittite empire, which have been carried on, with interruptions, since 1906.

The second phase began in 1925 with the commencement of American excavations at Yorgan Tepe (ancient Nuzi), 140 miles north of Baghdad, a provincial centre with Old Akkadian, Old Assyrian, and Middle Assyrian/Hurrian levels. There followed, among others, French excavations at Arslan Tash (ancient Hadatu; 1928), at Tall al-Aḥmar (ancient Til Barsib; 1929–31), and above all at Tall Ḥariri (ancient Mari; 1933 onward) and American excavations in the Diyālā region (east of Baghdad), at Tall al-Asmar (ancient Eshnunna), at Khafājī, and at other sites. Thus, excavation in Mesopotamia had moved away from the capital cities to include the “provinces.” Simultaneously, it expanded beyond the limits of Mesopotamia and Susiana and revealed outliers of “cuneiform civilization” on the Syrian coast at Ras Shamra (ancient Ugarit; France, 1929 onward) and on the Orontes of northern Syria at Al-ʿAṭshānāh (ancient Alalakh; England, 1937–39 and 1947–49), while, since 1954, Danish excavations on the islands of Bahrain and Faylakah, off the Tigris-Euphrates delta, have disclosed staging posts between Mesopotamia and the Indus Valley Civilization. Short-lived salvage operations have been undertaken at the site of the Assad Dam on the middle Euphrates (e.g., German excavations at Habūba al-Kabira, 1971–76). Italian excavations at Tall Mardikh (ancient Ebla; 1967 onward) have yielded spectacular results, including several thousand cuneiform tablets dating from the 24th century BC.

In its third phase, archaeological research in Mesopotamia and its neighbouring lands has probed back into prehistory and protohistory. The objective of these investigations, initiated by American archaeologists, was to trace as closely as possible the successive chronological stages in the progress of man from hunter-gatherer to settled farmer and, finally, to city dweller. These excavations are strongly influenced by the methods of the prehistorian, and the principal objective is no longer the search for texts and monuments. Apart from the American investigations, Iraq itself has taken part in this phase of the history of investigation, as has Japan since 1956 and the former Soviet Union from 1969 until the early 1990s.

Finally, the fourth category, which runs parallel with the first three phases, is represented by “surveys,” which do not concentrate on individual sites but attempt to define the relations between single settlements, their positioning along canals or rivers, or the distribution of central settlements and their satellites. Since shortages of time, money, and an adequate task force preclude the thorough investigation of large numbers of individual sites, the method employed is that of observing and collecting finds from the surface. Of these finds, the latest in date will give a rough termination date for the duration of the settlement, but, since objects from earlier, if not the earliest, levels work their way to the surface with a predictable degree of certainty or are exposed in rain gullies, an intensive search of the surface of the mound allows conclusions as to the total period of occupation with some degree of probability. If the individual periods of settlement are marked on superimposed maps, a very clear picture is obtained of the fluctuations in settlement patterns, of the changing proportions between large and small settlements, and of the equally changeable systems of riverbeds and irrigation canals—for, when points on the map lie in line, it is a legitimate assumption that they were once connected by watercourses.

During the four phases outlined, the objectives and methods of excavation have broadened and shifted. At first the chief aim was the recovery of valuable finds suitable for museums, but at the same time there was, from early on, considerable interest in the architecture of Mesopotamia, which has won for it the place it deserves in architectural history. Alongside philology, art history has also made great strides, building up a chronological framework by the combination of evidence from stratigraphic and stylistic criteria, particularly in pottery and cylinder seals. The discovery of graves and a variety of burial customs has thrown new light on the history of religion, stimulated by the interest of Old Testament studies. While pottery was previously collected for purely aesthetic motives or from the point of view of art history, attention has come to be paid increasingly to everyday wares, and greater insight into social and economic history is based on knowledge of the distribution and frequency of shapes and materials. The observation and investigation of animal bones and plant remains (pollen and seed analysis) have supplied invaluable information on the process of domestication, the conditions of animal husbandry, and the advances in agriculture. Such studies demand the cooperation of both zoologists and paleobotanists. In addition, microscopic analysis of the floors of excavated buildings may help to identify the functions of individual rooms.

The emergence of Mesopotamian civilization. *The Late Neolithic Period and the Chalcolithic Period.* Between about 10,000 BC and the genesis of large permanent settlements, the following stages of development are distinguishable, some of which run parallel: (1) the change to sedentary life, or the transition from continual or seasonal change of abode, characteristic of hunter-gatherers and the earliest cattle breeders, to life in one place over a period of several years or even permanently, (2) the transition from experimental plant cultivation to the deliberate and calculated farming of grains and leguminous plants, (3) the erection of houses and the associated “settlement” of the gods in temples, (4) the burial of the dead in cemeteries, (5) the invention of clay vessels, made at first by hand, then turned on the wheel and fired to ever greater degrees of hardness, at the same time receiving almost invariably

Advances
in Assyri-
ological
studies

Investi-
gations
after 1925

decoration of incised designs or painted patterns, (6) the development of specialized crafts and the distribution of labour, and (7) metal production (the first use of metal—copper—marks the transition from the Late Neolithic to the Chalcolithic Period).

These stages of development can only rarely be dated on the basis of a sequence of levels at one site alone. Instead, an important role is played by the comparison of different sites, starting with the assumption that what is simpler and technically less accomplished is older. In addition to this type of dating, which can be only relative, the radiocarbon, or carbon-14, method has proved to be an increasingly valuable tool since the 1950s. By this method the known rate of decay of the radioactive carbon isotope (carbon-14) in wood, horn, plant fibre, and bone allows the time that has elapsed since the “death” of the material under examination to be calculated. Although a plus/minus discrepancy of up to 200 years has to be allowed for, this is not such a great disadvantage in the case of material 6,000 to 10,000 years old. Even when skepticism is necessary because of the use of an inadequate sample, carbon-14 dates are still very welcome as confirmation of dates arrived at by other means. Moreover, radiocarbon ages can be converted to more precise dates through comparisons with data obtained by dendrochronology, a method of absolute age determination based on the analysis of the annual rings of trees.

The first agriculture, the domestication of animals, and the transition to sedentary life took place in regions in which animals that were easily domesticated, such as sheep, goats, cattle, and pigs, and the wild prototypes of grains and leguminous plants, such as wheat, barley, bitter vetch, pea, and lentil, were present. Such centres of dispersion may have been the valleys and grassy border regions of the mountains of Iran, Iraq, Anatolia, Syria, and Palestine, but they also could have been, say, the northern slopes of the Hindu Kush. As settled life, which caused a drop in infant mortality, led to the increase of the population, settlement spread out from these centres into the plains—although it must be remembered that this process, described as the Neolithic Revolution, in fact took thousands of years.

Representative of the first settlements on the borders of Mesopotamia are the adjacent sites of Zawi Chemi Shanidar and Shanidar itself, which lie northwest of Rawāndūz. They date from the transition from the 10th to the 9th millennium BC and are classified as prepottery. The finds included querns (primitive mills) for grinding grain (whether wild or cultivated is not known), the remains of huts about 13 feet in diameter, and a cemetery with grave goods. The presence of copper beads is evidence of acquaintance with metal, though not necessarily with the technique of working it into tools, and the presence of obsidian (volcanic glass) is indicative of the acquisition of nonindigenous raw materials by means of trade. The bones found testify that sheep were already domesticated at Zawi Chemi Shanidar.

At Karim Shahir, a site that cannot be accurately tied chronologically to Shanidar, clear proof was obtained both of the knowledge of grain cultivation, in the form of sickle blades showing sheen from use, and of the baking of clay, in the form of lightly fired clay figurines. Still in the hilly borders of Mesopotamia, a sequence of about 3,000 years can be followed at the site of Qal'at Jarmo, east of Kirkūk, some 150 miles north of Baghdad. The beginning of this settlement can be dated to about 6750 BC; excavations uncovered 12 archaeological levels of a regular village, consisting of about 20 to 25 houses built of packed clay, sometimes with stone foundations, and divided into several rooms. The finds included types of wheat (emmer and einkorn) and two-row barley, the bones of domesticated goats, sheep, and pigs, and obsidian tools, stone vessels, and, in the upper third of the levels, clay vessels with rough painted decorations, providing the first certain evidence for the manufacture of pottery. Jarmo must be roughly contemporary with the sites of Jericho (13 miles east of Jerusalem) and of Çatalhöyük in Anatolia (central Turkey). Those sites, with their walled settlements, seem to have achieved a much higher level of civilization, but

too much weight must not be placed on the comparison because no other sites in and around Mesopotamia confirm the picture deduced from Jarmo alone. Views on the earliest Neolithic in Iraq have undergone radical revisions in the light of discoveries made since the 1970s at Qermez Dere, Nemrik, and Maghzaliyah.

About 1,000 years later are two villages that are the earliest so far discovered in the plain of Mesopotamia: Ḥassūna, near Mosul, and Tall Ṣawwān, near Sāmarrā'. At Ḥassūna the pottery is more advanced, with incised and painted designs, but the decoration is still unsophisticated. One of the buildings found may be a shrine, judging from its unusual ground plan. Apart from emmer there occurs, as the result of mutation, six-row barley, which was later to become the chief grain crop of southern Mesopotamia. In the case of Tall Ṣawwān, it is significant that the settlement lay south of the boundary of rainfall agriculture; thus it must have been dependent on some form of artificial irrigation, even if this was no more than the drawing of water from the Tigris. This, therefore, gives a date after which the settlement of parts of southern Mesopotamia would have been feasible.

The emergence of cultures. For the next millennium, the 5th, it is customary to speak in terms of various “cultures” or “horizons,” distinguished in general by the pottery, which may be classed by its colour, shape, hardness, and, above all, by its decoration. The name of each horizon is derived either from the type site or from the place where the pottery was first found: Sāmarrā' on the Tigris, Tall Ḥalaf in the central Jazīrah, Ḥassūna Level V, Al-'Ubaid near Ur, and Ḥājj Muḥammad on the Euphrates, not far from As-Samāwah (some 150 miles south-southeast of Baghdad). Along with the improvement of tools, the first evidence for water transport (a model boat from the prehistoric cemetery at Eridu, in the extreme south of Mesopotamia, c. 4000 BC), and the development of terra-cottas, the most impressive sign of progress is the constantly accelerating advance in architecture. This can best be followed in the city of Eridu, which in historical times was the centre of the cult of the Sumerian god Enki.

Originally a small, single-roomed shrine, the temple in the Ubaid period consisted of a rectangular building, measuring 80 by 40 feet, that stood on an artificial terrace. It had an “offering table” and an “altar” against the short walls, aisles down each side, and a facade decorated with niches. This temple, standing on a terrace probably originally designed to protect the building from flooding, is usually considered the prototype of the characteristic religious structure of later Babylonia, the ziggurat. The temple at Eridu is in the very same place as that on which the Enki ziggurat stood in the time of the 3rd dynasty of Ur (c. 2112–c. 2004 BC), so the cult tradition must have existed on the same spot for at least 1,500 to 2,000 years before Ur III itself. Remarkable as this is, however, it is not justifiable to assume a continuous ethnic tradition. The flowering of architecture reached its peak with the great temples (or assembly halls?) of Uruk, built around the turn of the 4th to 3rd millennium BC (Uruk Levels VI to IV).

In extracting information as to the expression of mind and spirit during the six millennia preceding the invention of writing, it is necessary to take account of four major sources: decoration on pottery, the care of the dead, sculpture, and the designs on seals. There is, of course, no justification in assuming any association with ethnic groups.

The most varied of these means of expression is undoubtedly the decoration of pottery. It is hardly coincidental that, in regions in which writing had developed, high-quality painted pottery was no longer made. The motifs in decoration are either abstract and geometric or figured, although there is also a strong tendency to geometric stylization. An important question is the extent to which the presence of symbols, such as the bucranium (a sculptured ornament representing an ox skull), can be considered as expressions of specific religious ideas, such as a bull cult, and, indeed, how much the decoration was intended to convey meaning at all.

It is not known how ancient is the custom of burying the dead in graves nor whether its intention was to main-

The begin-
nings of
agriculture

Six-row
barley

Burial
customs

tain communication (by the cult of the dead) or to guard against the demonic power of the unburied dead left free to wander. A cemetery, or collection of burials associated with grave goods, is first attested at Zawi Chemi Shanidar. The presence of pots in the grave indicates that the bodily needs of the dead person were provided for, and the discovery of the skeleton of a dog and of a model boat in the cemetery at Eridu suggests that it was believed that the activities of life could be pursued in the afterlife.

The earliest sculpture takes the form of very crudely worked terra-cotta representations of women; the Ubaid Horizon, however, has figurines of both women and men, with very slender bodies, protruding features, arms akimbo, and the genitals accurately indicated, and also of women suckling children. It is uncertain whether it is correct to describe these statuettes as idols, whether the figures were cult objects, such as votive offerings, or whether they had a magical significance, such as fertility charms, or, indeed, what purpose they did fulfill.

Seals are first attested in the form of stamp seals at Tepe Gawra, north of Mosul. Geometric designs are found earlier than scenes with figures, such as men, animals, conflict between animals, copulation, or dance. Here again it is uncertain whether the scenes are intended to convey a deeper meaning. Nevertheless, unlike pottery, a seal has a direct relationship to a particular individual or group, for the seal identifies what it is used to seal (a vessel, sack, or other container) as the property or responsibility of a specific person. To that extent, seals represent the earliest pictorial representations of persons. The area of distribution of the stamp seal was northern Mesopotamia, Anatolia, and Iran. Southern Mesopotamia, on the other hand, was the home of the cylinder seal, which was either an independent invention or was derived from stamp seals engraved on two faces. The cylinder seal, with its greater surface area and more practical application, remained in use into the 1st millennium BC. Because of the continuous changes in the style of the seal designs, cylinder seals are among the most valuable of chronological indicators for archaeologists.

In general, the prehistory of Mesopotamia can only be described by listing and comparing human achievements, not by recounting the interaction of individuals or peoples. There is no basis for reconstructing the movements and migrations of peoples unless one is prepared to equate the spread of particular archaeological types with the extent of a particular population, the change of types with a change of population, or the appearance of new types with an immigration.

The only certain evidence for the movement of peoples beyond their own territorial limits is provided at first by material finds that are not indigenous. The discovery of obsidian and lapis lazuli at sites in Mesopotamia or in its neighbouring lands is evidence for the existence of trade, whether consisting of direct caravan trade or of a succession of intermediate stages.

Just as no ethnic identity is recognizable, so nothing is known of the social organization of prehistoric settlements. It is not possible to deduce anything of the "government" in a village nor of any supraregional connections that may have existed under the domination of one centre. Constructions that could only have been accomplished by the organization of workers in large numbers are first found in Uruk Levels VI to IV: the dimensions of these buildings suggest that they were intended for gatherings of hundreds of people. As for artificial irrigation, which was indispensable for agriculture in south Mesopotamia, the earliest form was probably not the irrigation canal. It is assumed that at first floodwater was dammed up to collect in basins, near which the fields were located. Canals, which led the water farther from the river, would have become necessary when the land in the vicinity of the river could no longer supply the needs of the population.

Mesopotamian protohistory. Attempts have been made by philologists to reach conclusions about the origin of the flowering of civilization in southern Mesopotamia by the analysis of Sumerian words. It has been thought possible to isolate an earlier, non-Sumerian substratum from the Sumerian vocabulary by assigning certain words on the

Beginnings
of artificial
irrigation

basis of their endings to either a Neolithic or a Chalcolithic language stratum. These attempts are based on the phonetic character of Sumerian at the beginning of the 2nd millennium BC, which is at least 1,000 years later than the invention of writing. Quite apart, therefore, from the fact that the structure of Sumerian words themselves is far from adequately investigated, the enormous gap in time casts grave doubt on the criteria used to distinguish between Sumerian and "pre-Sumerian" vocabulary.

The earliest peoples of Mesopotamia who can be identified from inscribed monuments and written tradition—people in the sense of speakers of a common language—are, apart from the Sumerians, Semitic peoples (Akkadians or pre-Akkadians) and Subarians (identical with, or near relatives of, the Hurrians, who appear in northern Mesopotamia around the end of the 3rd millennium BC). Their presence is known, but no definite statements about their past or possible routes of immigration are possible.

At the turn of the 4th to 3rd millennium BC, the long span of prehistory is over, and the threshold of the historical era is gained, captured by the existence of writing. Names, speech, and actions are fixed in a system that is composed of signs representing complete words or syllables. The signs may consist of realistic pictures, abbreviated representations, and perhaps symbols selected at random. Since clay is not well suited to the drawing of curved lines, a tendency to use straight lines rapidly gained ground. When the writer pressed the reed in harder at the beginning of a stroke, it made a triangular "head," and thus "wedges" were impressed into the clay. It is the Sumerians who are usually given the credit for the invention of this, the first system of writing in the Middle East. As far as they can be assigned to any language, the inscribed documents from before the dynasty of Akkad (c. 2334–c. 2154 BC) are almost exclusively in Sumerian. Moreover, the extension of the writing system to include the creation of syllabograms by the use of the sound of a logogram (sign representing a word), such as *gi*, "a reed stem," used to render the verb *gi*, "to return," can only be explained in terms of the Sumerian language. It is most probable, however, that Mesopotamia in the 4th millennium BC, just as in later times, was composed of many races. This makes it likely that, apart from the Sumerians, the interests and even initiatives of other language groups may have played their part in the formation of the writing system. Many scholars believe that certain clay objects or tokens that are found in prehistoric strata may have been used for some kind of primitive accounting. These tokens, some of which are incised and which have various forms, may thus be three-dimensional predecessors of writing.

Sumerian is an agglutinative language: prefixes and suffixes, which express various grammatical functions and relationships, are attached to a noun or verb root in a "chain." Attempts to identify Sumerian more closely by comparative methods have as yet been unsuccessful and will very probably remain so, as languages of a comparable type are known only from AD 500 (Georgian) or 1000 (Basque)—that is, 3,000 years later. Over so long a time, the rate of change in a language, particularly one that is not fixed in a written norm, is so great that one can no longer determine whether apparent similarity between words goes back to an original relationship or is merely fortuitous. Consequently, it is impossible to obtain any more accurate information as to the language group to which Sumerian may once have belonged.

The most important development in the course of the 4th millennium BC was the birth of the city. There were precursors, such as the unwalled prepottery settlement at Jericho of about 7000 BC, but the beginning of cities with a more permanent character came only later. There is no generally accepted definition of a city. In this context, it means a settlement that serves as a centre for smaller settlements, one that possesses one or more shrines of one or more major deities, has extensive granaries, and, finally, displays an advanced stage of specialization in the crafts.

The earliest cities of southern Mesopotamia, as far as their names are known, are Eridu, Uruk, Bad-tibira, Nippur, and Kish (35 miles south-southeast of Baghdad). The surveys of the American archaeologist Robert McCormick

Precursors
of writing

Adams and the German archaeologist Hans Nissen have shown how the relative size and number of the settlements gradually shifted: the number of small or very small settlements was reduced overall, whereas the number of larger places grew. The clearest sign of urbanization can be seen at Uruk, with the almost explosive increase in the size of the buildings. Uruk Levels VI to IV had rectangular buildings covering areas as large as 275 by 175 feet. These buildings are described as temples, since the ground plans are comparable to those of later buildings whose sacred character is beyond doubt, but other functions, such as assembly halls for noncultic purposes, cannot be excluded.

The major accomplishments of the period Uruk VI to IV, apart from the first inscribed tablets (Level IV B), are masterpieces of sculpture and of seal engraving and also of the form of wall decoration known as cone mosaics. Together with the everyday pottery of gray or red burnished ware, there is a very coarse type known as the beveled-rim bowl. These are vessels of standard size whose shape served as the original for the sign *sila*, meaning "litre." It is not too rash to deduce from the mass production of such standard vessels that they served for the issue of rations. This would have been the earliest instance of a system that remained typical of the southern Mesopotamian city for centuries: the maintenance of part of the population by allocations of food from the state.

Historians usually date the beginning of history, as opposed to prehistory and protohistory, from the first appearance of usable written sources. If this is taken to be the transition from the 4th to the 3rd millennium BC, it must be remembered that this applies only to part of Mesopotamia: the south, the Diyāla region, Susiana (with a later script of its own invented locally), and the district of the middle Euphrates, as well as Iran.

SUMERIAN CIVILIZATION

The Sumerians to the end of the Early Dynastic period. Despite the Sumerians' leading role, the historical role of other races should not be underestimated. While with prehistory only approximate dates can be offered, historical periods require a firm chronological framework, which, unfortunately, has not yet been established for the first half of the 3rd millennium BC. The basis for the chronology after about 1450 BC is provided by the data in the Assyrian and Babylonian king lists, which can often be checked by dated tablets and the Assyrian lists of eponyms (annual officials whose names served to identify each year). It is, however, still uncertain how much time separated the middle of the 15th century BC from the end of the 1st dynasty of Babylon, which is therefore variously dated to 1594 BC ("middle"), 1530 BC ("short"), or 1730 BC ("long" chronology). As a compromise, the middle chronology is used here. From 1594 BC several chronologically overlapping dynasties reach back to the beginning of the 3rd dynasty of Ur, about 2112 BC. From this point to the beginning of the dynasty of Akkad (c. 2334 BC) the interval can only be calculated to within 40 to 50 years, via the ruling houses of Lagash and the rather uncertain traditions regarding the succession of Gutian viceroys. With Ur-Nanshe (c. 2520 BC), the first king of the 1st dynasty of Lagash, there is a possible variation of 70 to 80 years, and earlier dates are a matter of mere guesswork: they depend upon factors of only limited relevance, such as the computation of occupation or destruction levels, the degree of development in the script (paleography), the character of the sculpture, pottery, and cylinder seals, and their correlation at different sites. In short, the chronology of the first half of the 3rd millennium is largely a matter for the intuition of the individual author. Carbon-14 dates are at present too few and far between to be given undue weight. Consequently, the turn of the 4th to 3rd millennium is to be accepted, with due caution and reservations, as the date of the flourishing of the archaic civilization of Uruk and of the invention of writing.

In Uruk and probably also in other cities of comparable size, the Sumerians led a city life that can be more or less reconstructed as follows: temples and residential districts; intensive agriculture, stock breeding, fishing, and date palm cultivation forming the four mainstays of the

economy; and highly specialized industries carried on by sculptors, seal engravers, smiths, carpenters, shipbuilders, potters, and workers of reeds and textiles. Part of the population was supported with rations from a central point of distribution, which relieved people of the necessity of providing their basic food themselves, in return for their work all day and every day, at least for most of the year. The cities kept up active trade with foreign lands.

That organized city life existed is demonstrated chiefly by the existence of inscribed tablets. The earliest tablets contain figures with the items they enumerate and measures with the items they measure, as well as personal names and, occasionally, probably professions. This shows the purely practical origins of writing in Mesopotamia: it began not as a means of magic or as a way for the ruler to record his achievements, for example, but as an aid to memory for an administration that was ever expanding its area of operations. The earliest examples of writing are very difficult to penetrate because of their extremely laconic formulation, which presupposes a knowledge of the context, and because of the still very imperfect rendering of the spoken word. Moreover, many of the archaic signs were pruned away after a short period of use and cannot be traced in the paleography of later periods, so that they cannot be identified.

One of the most important questions that has to be met when dealing with "organization" and "city life" is that of social structure and the form of government; however, it can be answered only with difficulty, and the use of evidence from later periods carries with it the danger of anachronisms. The Sumerian word for ruler par excellence is *lugal*, which etymologically means "big person." The first occurrence comes from Kish about 2700 BC, since an earlier instance from Uruk is uncertain because it could simply be intended as a personal name: "Monsieur Legrand." In Uruk the ruler's special title was *en*. In later periods this word (etymology unknown), which is also found in divine names such as Enlil and Enki, has a predominantly religious connotation that is translated, for want of a better designation, as "en-priest, en-priestess." *En*, as the ruler's title, is encountered in the traditional epics of the Sumerians (Gilgamesh is the "*en* of Kullab," a district of Uruk) and particularly in personal names, such as "The-*en*-has-abundance," "The-*en*-occupies-the-throne," and many others.

It has often been asked if the ruler of Uruk is to be recognized in artistic representations. A man feeding sheep with flowering branches, a prominent personality in seal designs, might thus represent the ruler or a priest in his capacity as administrator and protector of flocks. The same question may be posed in the case of a man who is depicted on a stela aiming an arrow at a lion. These questions are purely speculative, however: even if the "protector of flocks" were identical with the *en*, there is no ground for seeing in the ruler a person with a predominantly religious function.

Literary and other historical sources. The picture offered by the literary tradition of Mesopotamia is clearer but not necessarily historically relevant. The Sumerian king list has long been the greatest focus of interest. This is a literary composition, dating from Old Babylonian times, that describes kingship (*nam-lugal* in Sumerian) in Mesopotamia from primeval times to the end of the 1st dynasty of Isin. According to the theory—or rather the ideology—of this work, there was officially only one kingship in Mesopotamia, which was vested in one particular city at any one time; hence the change in dynasties brought with it the change of the seat of kingship:

Kish-Uruk-Ur-Awan-Kish-Hamazi-Uruk-Ur-
Adab-Mari-Kish-Akshak-Kish-Uruk-Akkad-
Uruk-Gutians-Uruk-Ur-Isin.

The king list gives as coming in succession several dynasties that now are known to have ruled simultaneously. It is a welcome aid to chronology and history, but, so far as the regnal years are concerned, it loses its value for the time before the dynasty of Akkad, for here the lengths of reign of single rulers are given as more than 100 and sometimes even several hundred years. One group of versions of the king list has adopted the tradition of the Sumerian Flood

Appear-
ance of
written
sources

Economic
and
administra-
tive texts

The
Sumerian
king list

story, according to which Kish was the first seat of kingship after the Flood, whereas five dynasties of primeval kings ruled before the Flood in Eridu, Bad-tibira, Larak, Sippar, and Shuruppak. These kings all allegedly ruled for multiples of 3,600 years (the maximum being 64,800 or, according to one variant, 72,000 years). The tradition of the Sumerian king list is still echoed in Berossus.

It is also instructive to observe what the Sumerian king list does not mention. The list lacks all mention of a dynasty as important as the 1st dynasty of Lagash (from King Ur-Nanshe to Urukagina) and appears to retain no memory of the archaic florescence of Uruk at the beginning of the 3rd millennium BC.

Besides the peaceful pursuits reflected in art and writing, the art also provides the first information about violent contacts: cylinder seals of the Uruk Level IV depict fettered men lying or squatting on the ground, being beaten with sticks or otherwise maltreated by standing figures. They may represent the execution of prisoners of war. It is not known from where these captives came or what form "war" would have taken or how early organized battles were fought. Nevertheless, this does give the first, albeit indirect, evidence for the wars that are henceforth one of the most characteristic phenomena in the history of Mesopotamia.

Just as with the rule of man over man, with the rule of higher powers over man it is difficult to make any statements about the earliest attested forms of religion or about the deities and their names without running the risk of anachronism. Excluding prehistoric figurines, which provide no evidence for determining whether men or anthropomorphic gods are represented, the earliest testimony is supplied by certain symbols that later became the cuneiform signs for gods' names: the "gatepost with streamers" for Inanna, goddess of love and war, and the "ringed post" for the moon god Nanna. A scene on a cylinder seal—a shrine with an Inanna symbol and a "man" in a boat—could be an abbreviated illustration of a procession of gods or of a cultic journey by ship. The constant association of the "gatepost with streamers" with sheep and of the "ringed post" with cattle may possibly reflect the area of responsibility of each deity. The Sumerologist Thorkild Jacobsen sees in the pantheon a reflex of the various economies and modes of life in ancient Mesopotamia: fishermen and marsh dwellers, date palm cultivators, cowherds, shepherds, and farmers all have their special groups of gods.

Both Sumerian and non-Sumerian languages can be detected in the divine names and place-names. Since the pronunciation of the names is known only from 2000 BC or later, conclusions about their linguistic affinity are not without problems. Several names, for example, have been reinterpreted in Sumerian by popular etymology. It would be particularly important to isolate the Subarian components (related to Hurrian), whose significance was probably greater than has hitherto been assumed. For the south Mesopotamian city HA.A (the noncommittal transliteration of the signs) there is a pronunciation gloss "shubari," and non-Sumerian incantations are known in the language of HA.A that have turned out to be "Subarian."

There have always been in Mesopotamia speakers of Semitic languages (which belong to the Afro-Asiatic group and also include ancient Egyptian, Berber, and various African languages). This element is easier to detect in ancient Mesopotamia, but whether people began to participate in city civilization in the 4th millennium BC or only during the 3rd is unknown. Over the last 4,000 years, Semites (Amorites, Canaanites, Aramaeans, and Arabs) have been partly nomadic, ranging the Arabian fringes of the Fertile Crescent, and partly settled; and the transition to settled life can be observed in a constant, though uneven, rhythm. There are, therefore, good grounds for assuming that the Akkadians (and other pre-Akkadian Semitic tribes not known by name) also originally led a nomadic life to a greater or lesser degree. Nevertheless, they can only have been herders of domesticated sheep and goats, which require changes of pasturage according to the time of year and can never stray more than a day's march from the watering places. The traditional nomadic

life of the Bedouin makes its appearance only with the domestication of the camel at the turn of the 2nd to 1st millennium BC.

The question arises as to how quickly writing spread and by whom it was adopted in about 3000 BC or shortly thereafter. At Kish, in northern Babylonia, almost 120 miles northwest of Uruk, a stone tablet has been found with the same repertoire of archaic signs as those found at Uruk itself. This fact demonstrates that intellectual contacts existed between northern and southern Babylonia. The dispersion of writing in an unaltered form presupposes the existence of schools in various cities that worked according to the same principles and adhered to one and the same canonical repertoire of signs. It would be wrong to assume that Sumerian was spoken throughout the area in which writing had been adopted. Moreover, the use of cuneiform for a non-Sumerian language can be demonstrated with certainty from the 27th century BC.

First historical personalities. The specifically political events in Mesopotamia after the flourishing of the archaic culture of Uruk cannot be pinpointed. Not until about 2700 BC does the first historical personality appear—historical because his name, Enmebaragesi (Me-baragesi), was preserved in later tradition. It has been assumed, although the exact circumstances cannot be reconstructed, that there was a rather abrupt end to the high culture of Uruk Level IV. The reason for the assumption is a marked break in both artistic and architectural traditions: entirely new styles of cylinder seals were introduced; the great temples (if in fact they were temples) were abandoned, flouting the rule of a continuous tradition on religious sites, and on a new site a shrine was built on a terrace, which was to constitute the lowest stage of the later Eanna ziggurat. On the other hand, since the writing system developed organically and was continually refined by innovations and progressive reforms, it would be overhasty to assume a revolutionary change in the population.

In the quarter or third of a millennium between Uruk Level IV and Enmebaragesi, southern Mesopotamia became studded with a complex pattern of cities, many of which were the centres of small independent city-states, to judge from the situation in about the middle of the millennium. In these cities, the central point was the temple, sometimes encircled by an oval boundary wall (hence the term temple oval); but nonreligious buildings, such as palaces serving as the residences of the rulers, could also function as centres.

Enmebaragesi, king of Kish, is the oldest Mesopotamian ruler from whom there are authentic inscriptions. These are vase fragments, one of them found in the temple oval of Khafajah (Khafāji). In the Sumerian king list, Enmebaragesi is listed as the penultimate king of the 1st dynasty of Kish; a Sumerian poem, "Gilgamesh and Agga of Kish," describes the siege of Uruk by Agga, son of Enmebaragesi. The discovery of the original vase inscriptions was of great significance because it enabled scholars to ask with somewhat more justification whether Gilgamesh, the heroic figure of Mesopotamia who has entered world literature, was actually a historical personage. The indirect synchronism notwithstanding, the possibility exists that even remote antiquity knew its "Ninus" and its "Semi-ramis," figures onto which a rapidly fading historical memory projected all manner of deeds and adventures. Thus, though the historical tradition of the early 2nd millennium believes Gilgamesh to have been the builder of the oldest city wall of Uruk, such may not have been the case. The palace archives of Shuruppak (modern Tall Fa'rah, 125 miles southeast of Baghdad), dating presumably from shortly after 2600, contain a long list of divinities, including Gilgamesh and his father Lugalbanda. More recent tradition, on the other hand, knows Gilgamesh as judge of the nether world. However that may be, an armed conflict between two Mesopotamian cities such as Uruk and Kish would hardly have been unusual in a country whose energies were consumed, almost without interruption from 2500 to 1500 BC, by clashes between various separatist forces. The great "empires," after all, formed the exception, not the rule.

Emergent city-states. Kish must have played a major

The spread of writing

Early religious beliefs

Gilgamesh

role almost from the beginning. After 2500, southern Babylonian rulers, such as Mesannepada of Ur and Eannatum of Lagash, frequently called themselves king of Kish when laying claim to sovereignty over northern Babylonia. This does not agree with some recent histories in which Kish is represented as an archaic "empire." It is more likely to have figured as representative of the north, calling forth perhaps the same geographic connotation later evoked by "the land of Akkad."

Although the corpus of inscriptions grows richer both in geographic distribution and in point of chronology in the 27th and increasingly so in the 26th century, it is still impossible to find the key to a plausible historical account, and history cannot be written solely on the basis of archaeological findings. Unless clarified by written documents, such findings contain as many riddles as they seem to offer solutions. This applies even to as spectacular a discovery as that of the royal tombs of Ur with their hecatombs (large-scale sacrifices) of retainers who followed their king and queen to the grave, not to mention the elaborate funerary appointments with their inventory of tombs. It is only from about 2520 to the beginnings of the dynasty of Akkad that history can be written within a framework, with the aid of reports about the city-state of Lagash and its capital of Girsu and its relations with its neighbour and rival, Umma.

Sources for this are, on the one hand, an extensive corpus of inscriptions relating to nine rulers, telling of the buildings they constructed, of their institutions and wars, and, in the case of UruKagina, of their "social" measures. On the other hand, there is the archive of some 1,200 tablets—insofar as these have been published—from the temple of Baba, the city goddess of Girsu, from the period of Lugalanda and UruKagina (first half of the 24th century). For generations, Lagash and Umma contested the possession and agricultural usufruct of the fertile region of Gu'edena. To begin with, some two generations before Ur-Nanshe, Mesilim (another "king of Kish") had intervened as arbiter and possibly overlord in dictating to both states the course of the boundary between them, but this was not effective for long. After a prolonged struggle, Eannatum forced the ruler of Umma, by having him take an involved oath to six divinities, to desist from crossing the old border, a dike. The text that relates this event, with considerable literary elaboration, is found on the Stele of Vultures. These battles, favouring now one side, now the other, continued under Eannatum's successors, in particular Entemena, until, under UruKagina, great damage was done to the land of Lagash and to its holy places. The enemy, Lugalzagesi, was vanquished in turn by Sargon of Akkad. The rivalry between Lagash and Umma, however, must not be considered in isolation. Other cities, too, are occasionally named as enemies, and the whole situation resembles the pattern of changing coalitions and short-lived alliances between cities of more recent times. Kish, Umma, and distant Mari on the middle Euphrates are listed together on one occasion as early as the time of Eannatum. For the most part, these battles were fought by infantry, although mention is also made of war chariots drawn by onagers (wild asses).

The lords of Lagash rarely fail to call themselves by the title of *ensi*, of as yet undetermined derivation; "city ruler," or "prince," are only approximate translations. Only seldom do they call themselves *lugal*, or "king," the title given the rulers of Umma in their own inscriptions. In all likelihood, these were local titles that were eventually converted, beginning perhaps with the kings of Akkad, into a hierarchy in which the *lugal* took precedence over the *ensi*.

Territorial states. More difficult than describing its external relations is the task of shedding light on the internal structure of a state like Lagash. For the first time, a state consisting of more than a city with its surrounding territory came into being, because aggressively minded rulers had managed to extend that territory until it comprised not only Girsu, the capital, and the cities of Lagash and Nina (Zurghul) but also many smaller localities and even a seaport, Guabba. Yet it is not clear to what extent the conquered regions were also integrated administratively. On one occasion UruKagina used the formula "from the

limits of Ningirsu [that is, the city god of Girsu] to the sea," having in mind a distance of up to 125 miles. It would be unwise to harbour any exaggerated notion of well-organized states exceeding that size.

For many years, scholarly views were conditioned by the concept of the Sumerian temple city, which was used to convey the idea of an organism whose ruler, as representative of his god, theoretically owned all land, privately held agricultural land being a rare exception. The concept of the temple city had its origin partly in the overinterpretation of a passage in the so-called reform texts of UruKagina, that states "on the field of the *ensi* [or his wife and the crown prince], the city god Ningirsu [or the city goddess Baba and the divine couple's son]" had been "reinstated as owners." On the other hand, the statements in the archives of the temple of Baba in Girsu, dating from Lugalanda and UruKagina, were held to be altogether representative. Here is a system of administration, directed by the *ensi*'s spouse or by a *sangu* (head steward of a temple), in which every economic process, including commerce, stands in a direct relationship to the temple: agriculture, vegetable gardening, tree farming, cattle raising and the processing of animal products, fishing, and the payment in merchandise of workers and employees.

The conclusion from this analogy proved to be dangerous because the archives of the temple of Baba provide information about only a portion of the total temple administration and that portion, furthermore, is limited in time. Understandably enough, the private sector, which of course was not controlled by the temple, is scarcely mentioned at all in these archives. The existence of such a sector is nevertheless documented by bills of sale for land purchases of the pre-Sargonic period and from various localities. Written in Sumerian as well as in Akkadian, they prove the existence of private land ownership or, in the opinion of some scholars, of lands predominantly held as undivided family property. Although a substantial part of the population was forced to work for the temple and drew its pay and board from it, it is not yet known whether it was year-round work.

It is probable, if unfortunate, that there will never exist a detailed and numerically accurate picture of the demographic structure of a Sumerian city. It is assumed that in the oldest cities the government was in a position to summon sections of the populace for the performance of public works. The construction of monumental buildings or the excavation of long and deep canals could be carried out only by means of such a levy. The large-scale employment of indentured persons and of slaves is of no concern in this context. Evidence of male slavery is fairly rare before Ur III, and even in Ur III and in the Old Babylonian period slave labour was never an economically relevant factor. It was different with female slaves. According to one document, the temple of Baba employed 188 such women; the temple of the goddess Nanshe employed 180, chiefly in grinding flour and in the textile industry, and this continued to be the case in later times. For accuracy's sake it should be added that the terms male slave and female slave are used here in the significance they possessed about 2000 and later, designating persons in bondage who were bought and sold and who could not acquire personal property through their labour. A distinction is made between captured slaves (prisoners of war and kidnapped persons) and others who had been sold.

In one inscription, Entemena of Lagash boasts of having "allowed the sons of Uruk, Larsa, and Bad-tibira to return to their mothers" and of having "restored them into the hands" of the respective city god or goddess. Read in the light of similar but more explicit statements of later date, this laconic formula represents the oldest known evidence of the fact that the ruler occasionally endeavoured to mitigate social injustices by means of a decree. Such decrees might refer to the suspension or complete cancellation of debts or to exemption from public works. Whereas a set of inscriptions of the last ruler from the 1st dynasty of Lagash, UruKagina, has long been considered a prime document of social reform in the 3rd millennium, the designation "reform texts" is only partly justified. Reading between the lines, it is possible to discern that tensions had

Temple
administra-
tion

Intercity
rivalry
between
Lagash and
Umma

Inscrip-
tions of
Uru-
Kagina

arisen between the "palace"—the ruler's residence with its annex, administrative staff, and landed properties—and the "clergy"—that is, the stewards and priests of the temples. In seeming defiance of his own interests, UruKagina, who in contrast to practically all of his predecessors lists no genealogy and has therefore been suspected of having been a usurper, defends the clergy, whose plight he describes somewhat tearfully.

If the foregoing passage about restoring the *ensi's* fields to the divinity is interpreted carefully, it would follow that the situation of the temple was ameliorated and that palace lands were assigned to the priests. Along with these measures, which resemble the policies of a newcomer forced to lean on a specific party, are found others that do merit the designation of "measures taken toward the alleviation of social injustices"—for instance, the granting of delays in the payment of debts or their outright cancellation and the setting up of prohibitions to keep the economically or socially more powerful from forcing his inferior to sell his house, his ass's foal, and the like. Besides this, there were tariff regulations, such as newly established fees for weddings and burials, as well as the precise regulation of the food rations of garden workers.

These conditions, described on the basis of source materials from Girsu, may well have been paralleled elsewhere, but it is equally possible that other archives, yet to be found in other cities of pre-Sargonic southern Mesopotamia, may furnish entirely new historical aspects. At any rate, it is wiser to proceed cautiously, keeping to analysis and evaluation of the available material rather than making generalizations.

This, then, is the horizon of Mesopotamia shortly before the rise of the Akkadian empire. In Mari, writing was introduced at the latest about the mid-26th century BC, and from that time this city, situated on the middle Euphrates, forms an important centre of cuneiform civilization, especially in regard to its Semitic component. Ebla (and probably many other sites in ancient Syria) profited from the influence of Mari scribal schools. Reaching out across the Diyālā region and the Persian Gulf, Mesopotamian influences extended to Iran, where Susa is mentioned along with Elam and other, not yet localized, towns. In the west the Amanus Mountains were known, and under Lugalzagesi the "upper sea"—in other words, the Mediterranean—is mentioned for the first time. To the east the inscriptions of Ur-Nanshe of Lagash name the isle of Dilmun (modern Bahrain), which may have been even then a transshipment point for trade with the Oman coast and the Indus region, the Magan and Meluhha of more recent texts. Trade with Anatolia and Afghanistan was nothing new in the 3rd millennium, even if these regions are not yet listed by their names. It was the task of the Akkadian dynasty to unite within these boundaries a territory that transcended the dimensions of a state of the type represented by Lagash.

Sumer and Akkad from 2350 to 2000 BC. There are several reasons for taking the year 2350 as a turning point in the history of Mesopotamia. For the first time, an empire arose on Mesopotamian soil. The driving force of that empire was the Akkadians, so called after the city of Akkad, which Sargon chose for his capital (it has not yet been identified but was presumably located on the Euphrates between Sippar and Kish). The name Akkad became synonymous with a population group that stood side by side with the Sumerians. Southern Mesopotamia became known as the "land of Sumer and Akkad"; Akkadian became the name of a language; and the arts rose to new heights. However, even this turning point was not the first time the Akkadians had emerged in history. Semites—whether Akkadians or a Semitic language group that had settled before them—may have had a part in the urbanization that took place at the end of the 4th millennium. The earliest Akkadian names and words occur in written sources of the 27th century. The names of several Akkadian scribes are found in the archives of Tall Abū Ṣalābikh, near Nippur in central Babylonia, synchronous with those of Shuruppak (shortly after 2600). The Sumerian king list places the 1st dynasty of Kish, together with a series of kings bearing Akkadian names, immediately after

the Flood. In Mari the Akkadian language was probably written from the very beginning. Thus, the founders of the dynasty of Akkad were presumably members of a people who had been familiar for centuries with Mesopotamian culture in all its forms.

Sargon's reign. According to the Sumerian king list, the first five rulers of Akkad (Sargon, Rimush, Manishtusu, Naram-Sin, and Shar-kali-sharri) ruled for a total of 142 years; Sargon alone ruled for 56. Although these figures cannot be checked, they are probably trustworthy, because the king list for Ur III, even if 250 years later, did transmit dates that proved to be accurate.

As stated in an annotation to his name in the king list, Sargon started out as a cupbearer to King Ur-Zababa of Kish. There is an Akkadian legend about Sargon, describing how he was exposed after birth, brought up by a gardener, and later beloved by the goddess Ishtar. Nevertheless, there are no historical data about his career. Yet it is feasible to assume that in his case a high court office served as springboard for a dynasty of his own. The original inscriptions of the kings of Akkad that have come down to posterity are brief, and their geographic distribution generally is more informative than is their content. The main sources for Sargon's reign, with its high points and catastrophes, are copies made by Old Babylonian scribes in Nippur from the very extensive originals that presumably had been kept there. They are in part Akkadian, in part bilingual Sumerian-Akkadian texts. According to these texts, Sargon fought against the Sumerian cities of southern Babylonia, threw down city walls, took prisoner 50 *ensis*, and "cleansed his weapons in the sea." He is also said to have captured Lugalzagesi of Uruk, the former ruler of Umma, who had vigorously attacked UruKagina in Lagash, forcing his neck under a yoke and leading him thus to the gate of the god Enlil at Nippur. "Citizens of Akkad" filled the offices of *ensi* from the "nether sea" (the Persian Gulf) upward, which was perhaps a device used by Sargon to further his dynastic aims. Aside from the 34 battles fought in the south, Sargon also tells of conquests in northern Mesopotamia: Mari, Tuttul on the Balikh, where he venerated the god Dagan (Dagon), Ebla (Tall Mardikh in Syria), the "cedar forest" (Amanus or Lebanon), and the "silver mountains"; battles in Elam and the foothills of the Zagros are mentioned. Sargon also relates that ships from Meluhha (Indus region), Magan (possibly the coast of Oman), and Dilmun (Bahrain) made fast in the port of Akkad.

Impressive as they are at first sight, these reports have only a limited value because they cannot be arranged chronologically, and it is not known whether Sargon built a large empire. Akkadian tradition itself saw it in this light, however, and a learned treatise of the late 8th or the 7th century lists no fewer than 65 cities and lands belonging to that empire. Yet, even if Magan and Kapturu (Crete) are given as the eastern and western limits of the conquered territories, it is impossible to transpose this to the 3rd millennium.

Sargon appointed one of his daughters priestess of the moon god in Ur. She took the name of Enheduanna and was succeeded in the same office by Enmenanna, a daughter of Naram-Sin. Enheduanna must have been a very gifted woman; two Sumerian hymns by her have been preserved, and she is also said to have been instrumental in starting a collection of songs dedicated to the temples of Babylonia.

Sargon died at a very old age. The inscriptions, also preserved only in copies, of his son Rimush are full of reports about battles fought in Sumer and Iran, just as if there had never been a Sargonic empire. It is not known in detail how rigorously Akkad wished to control the cities to the south and how much freedom had been left to them; but they presumably clung tenaciously to their inherited local autonomy. From a practical point of view, it was probably in any case impossible to organize an empire that would embrace all Mesopotamia.

Since the reports (*i.e.*, copies of inscriptions) left by Manishtusu, Naram-Sin, and Shar-kali-sharri speak time and again of rebellions and victorious battles and since Rimush, Manishtusu, and Shar-kali-sharri are themselves

Sargon's
conquests

Rise of the
Akkadian
empire

said to have died violent deaths, the problem of what remained of Akkad's greatness obtrudes. Wars and disturbances, the victory of one and the defeat of another, and even regicide constitute only some of the aspects suggested to us by the sources. Whenever they extended beyond the immediate Babylonian neighbourhood, the military campaigns of the Akkadian kings were dictated primarily by trade interests instead of being intended to serve the conquest and safeguarding of an empire. Akkad, or more precisely the king, needed merchandise, money, and gold in order to finance wars, buildings, and the system of administration that he had instituted.

On the other hand, the original inscriptions that have been found so far of a king like Naram-Sin are scattered at sites covering a distance of some 620 miles as the crow flies, following the Tigris downriver: Diyarbakır on the upper Tigris, Nineveh, Tall Birāk (Tell Brak) on the upper Khābūr River (which had an Akkadian fortress and garrison), Susa in Elam, as well as Marad, Puzrish-Dagan, Adab (Bismāyah), Nippur, Ur, and Girsu in Babylonia. Even if all this was not part of an empire, it surely constituted an impressive sphere of influence.

Also to be considered are other facts that weigh more heavily than high-sounding reports of victories that cannot be verified. After the first kings of the dynasty had borne the title of king of Kish, Naram-Sin assumed the title "king of the four quarters of the earth"—that is, of the universe. As if he were in fact divine, he also had his name written with the cuneiform sign "god," the divine determinative that was customarily used in front of the names of gods; furthermore, he assumed the title of "god of Akkad." It is legitimate to ask whether the concept of deification may be used in the sense of elevation to a rank equal to that of the gods. At the very least it must be acknowledged that, in relation to his city and his subjects, the king saw himself in the role played by the local divinity as protector of the city and guarantor of its well-being. In contemporary judicial documents from Nippur, the oath is often taken "by Naram-Sin," with a formula identical with that used in swearing by a divinity. Documents from Girsu contain Akkadian date formulas of the type "in the year in which Naram-Sin laid the foundations of the Enlil temple at Nippur and of the Inanna temple at Zabalam." As evidenced by the dating procedures customary in Ur III and in the Old Babylonian period, the use of such formulas presupposes that the respective city acknowledged as its overlord the ruler whose name is invoked.

Ascendancy of Akkad. Under Akkad, the Akkadian language acquired a literary prestige that made it the equal of Sumerian. Under the influence, perhaps, of an Akkadian garrison at Susa, it spread beyond the borders of Mesopotamia. After having employed for several centuries an indigenous script patterned after cuneiform writing, Elam adopted Mesopotamian script during the Akkadian period and with a few exceptions used it even when writing in Elamite rather than Sumerian or Akkadian. The so-called Old Akkadian manner of writing is extraordinarily appealing from the aesthetic point of view; as late as the Old Babylonian era it served as a model for monumental inscriptions. Similarly, the plastic and graphic arts, especially sculpture in the round, relief work, and cylinder seals, reached a high point of perfection.

Thus the reign of the five kings of Akkad may be considered one of the most productive periods of Mesopotamian history. Although separatist forces opposed all unifying tendencies, Akkad brought about a broadening of political horizons and dimensions. The period of Akkad fascinated historiographers as did few other eras. Having contributed its share to the storehouse of legend, it has never disappeared from memory. With phrases such as "There will come a king of the four quarters of the earth," liver omens (soothsaying done by analyzing the shape of a sheep's liver) of the Old Babylonian period express the yearning for unity at a time when Babylonia had once again disintegrated into a dozen or more small states.

The end of the dynasty. Of the kings after Shar-kali-sharri (c. 2217–c. 2193), only the names and a few brief inscriptions have survived. Quarrels arose over the succession, and the dynasty went under, although modern

scholars know as little about the individual stages of this decline as about the rise of Akkad. Two factors contributed to its downfall: the invasion of the nomadic Amurrus (Amorites), called Martu by the Sumerians, from the northwest, and the infiltration of the Gutians, who came, apparently, from the region between the Tigris and the Zagros Mountains to the east. This argument, however, may be a vicious circle, as these invasions were provoked and facilitated by the very weakness of Akkad. In Ur III the Amorites, in part already sedentary, formed one ethnic component along with Sumerians and Akkadians. The Gutians, on the other hand, played only a temporary role, even if the memory of a Gutian dynasty persisted until the end of the 17th century BC. As a matter of fact, the wholly negative opinion that even some modern historians have of the Gutians is based solely on a few stereotyped statements by the Sumerians and Akkadians, especially on the victory inscription of Utu-hegal of Uruk (c. 2116–c. 2110). While Old Babylonian sources give the region between the Tigris and the Zagros Mountains as the home of the Gutians, these people probably also lived on the middle Euphrates during the 3rd millennium. According to the Sumerian king list, the Gutians held the "kingship" in southern Mesopotamia for about 100 years. It has long been recognized that there is no question of a whole century of undivided Gutian rule and that some 50 years of this rule coincided with the final half century of Akkad. From this period there has also been preserved a record of a "Gutian interpreter." As it is altogether doubtful whether the Gutians had made any city of southern Mesopotamia their "capital" instead of controlling Babylonia more or less informally from outside, scholars cautiously refer to "viceroys" of this people. The Gutians have left no material records, and the original inscriptions about them are so scanty that no binding statements about them are possible.

The Gutians' influence probably did not extend beyond Umma. The neighbouring state of Lagash enjoyed a century of complete independence, between Shar-kali-sharri and the beginning of Ur III, during which time it showed expansionist tendencies and had widely ranging trade connections. Of the *ensi* Gudea, a contemporary of Ur-Nammu of Ur III, there are extant writings, exclusively Sumerian in language, which are of inestimable value. He had the time, power, and means to carry out an extensive program of temple construction during his reign, and in a hymn divided into two parts and preserved in two clay cylinders 12 inches (30 centimetres) high he describes explicitly the reconstruction of Eninnu, the temple of the god Ningirsu. Comprising 1,363 lines, the text is second in length only to Eannatum's Stele of Vultures among the literary works of the Sumerians up to that time. While Gudea forges a link, in his literary style, with his country's pre-Sargonic period, his work also bears the unmistakable stamp of the period of Akkad. Thus, the regions that furnish him building materials reflect the geographic horizon of the empire of Akkad, and the *ensi's* title "god of his city" recalls the "god of Akkad" (Naram-Sin). The building hymn contains interesting particulars about the work force deployed. "Levies" were organized in various parts of the country, and the city of Girsu itself "followed the *ensi* as though it were a single man." Unfortunately lacking are synchronous administrative archives of sufficient length to provide less summarily compiled information about the social structure of Lagash at the beginning of the 3rd dynasty of Ur. After the great pre-Sargonic archives of the Baba temple at Girsu, only the various administrative archives of the kings of Ur III give a closer look at the functioning of a Mesopotamian state.

The 3rd dynasty of Ur. Utu-hegal of Uruk is given credit for having overthrown Gutian rule by vanquishing their king Tiriqan along with two generals. Utu-hegal calls himself lord of the four quarters of the earth in an inscription, but this title, adopted from Akkad, is more likely to signify political aspiration than actual rule. Utu-hegal was a brother of the Ur-Nammu who founded the 3rd dynasty of Ur ("3rd" because it is the third time that Ur is listed in the Sumerian king list). Under Ur-Nammu and his successors Shulgi, Amar-Su'ena, Shu-Sin, and Ibbi-Sin, this

The
Gutians

Divine
kingship

Gudea of
Lagash

dynasty lasted for a century (c. 2112–c. 2004). Ur-Nammu was at first “governor” of the city of Ur under Utu-hegal. How he became king is not known, but there may well be some parallels between his rise and the career of Ishbi-Erra of Isin or, indeed, that of Sargon. By eliminating the state of Lagash, Ur-Nammu caused the coveted overseas trade (Dilmun, Magan, and Meluhha) to flow through Ur. As evidenced by a new royal title that he was the first to bear—that of “king of Sumer and Akkad”—he had built up a state that comprised at least the southern part of Mesopotamia. Like all great rulers, he built much, including the very impressive ziggurats of Ur and Uruk, which acquired their final monumental dimensions in his reign.

Assyriologists have given the name of Code of Ur-Nammu to a literary monument that is the oldest known example of a genre extending through the Code of Lipit-Ishtar in Sumerian to the Code of Hammurabi, written in Akkadian. (Some scholars have attributed it to Ur-Nammu’s son Shulgi.) It is a collection of sentences or verdicts mostly following the pattern of “If A [assumption], it follows that B [legal consequence].” The collection is framed by a prologue and an epilogue. The original was most likely a stela, but all that is known of the Code of Ur-Nammu so far are Old Babylonian copies. The term code as used here is conventional terminology and should not give the impression of any kind of “codified” law; furthermore, the content of the Code of Ur-Nammu is not yet completely known. It deals, among other things, with adultery by a married woman, the defloration of someone else’s female slave, divorce, false accusation, the escape of slaves, bodily injury, and the granting of security, as well as with legal cases arising from agriculture and irrigation.

Before its catastrophic end under Ibbi-Sin, the state of Ur III does not seem to have suffered setbacks and rebellions as grievous as those experienced by Akkad. There are no clear indications pointing to inner unrest, although it must be remembered that the first 20 years of Shulgi’s reign are still hidden in darkness. However, from that point on until the beginning of Ibbi-Sin’s reign, or for a period of 50 years at least, the sources give the impression of peace enjoyed by a country that lived undisturbed by encroachments from abroad. Some expeditions were sent into foreign lands, to the region bordering on the Zagros, to what later became Assyria, and to the vicinity of Elam, in order to secure the importation of raw materials, in a fashion reminiscent of Akkad. Force seems to have been employed only as a last resort, and every attempt was made to bring about peaceful conditions on the other side of the border through the dispatch of embassies or the establishment of family bonds—for example, by marrying the king’s daughters to foreign rulers.

Shulgi, too, called himself king of the four quarters of the earth. Although he resided in Ur, another important centre was in Nippur, whence—according to the prevailing ideology—Enlil, the chief god in the Sumerian state pantheon, had bestowed on Shulgi the royal dignity. Shulgi and his successors enjoyed divine honours, as Naram-Sin of Akkad had before them; by now, however, the process of deification had taken on clearer outlines in that sacrifices were offered and chapels built to the king and his throne, while the royal determinative turned up in personal names. Along with an Utu-hegal (“The Sun God Is Exuberance”) there appears a Shulgi-hegal (“Shulgi Is Exuberance”), and so forth.

Administration. The highest official of the state was the *sukkal-mah*, literally “supreme courier,” whose position may be described as “(state) chancellor.” The empire was divided into some 40 provinces ruled by as many *ensis*, who, despite their far-reaching authority (civil administration and judicial powers), were no longer autonomous, even if only indirectly, although the office was occasionally handed down from father to son. They could not enter into alliances or wage wars on their own. The *ensis* were appointed by the king and could probably also be transferred by him to other provinces. Each of these provinces was obliged to pay a yearly tribute, the amount of which was negotiated by emissaries. Of special significance in this was a system called *bala*, “cycle” or “rotation,” in which the *ensis* of the southern provinces took part; among other

things, they had to keep the state stockyards supplied with sacrificial animals. Although the “province” often corresponded to a former city-state, many others were no doubt newly established. The so-called land-register text of Ur-Nammu describes four such provinces north of Nippur, giving the precise boundaries and ending in each case with the statement, “King Ur-Nammu has confirmed the field of the god XX for the god XX.” In some cities, notably in Uruk, Mari, or Dēr (near Badrah, Iraq), the administration was in the hands of a *šakkana*, a man whose title is rendered partly by “governor” and partly by “general.”

The available histories are practically unanimous in seeing in Ur III a strongly centralized state marked by the king’s position as absolute ruler. Nevertheless, some caution is indicated. For one thing, the need to deal as carefully as possible with the *ensis* must not be underestimated. A further question arises from the borders between and relative extent of the “public” and the “private” sector; the latter’s importance may have been underrated as well. What is meant by “private” sector is a population group with land of its own and with revenues not directly granted by a temple or a “palace,” such as by the king’s or an *ensi*’s household. The traditional picture is derived from the sources, the state archives of Puzrish-Dagan, a gigantic “stockyard” situated outside the gates of Nippur, which supplied the city’s temples with sacrificial animals but inevitably also comprised a major wool and leather industry; other such archives are those of Umma, Girsu, Nippur, and Ur. All these activities were overseen by a finely honed bureaucracy that stressed the use of official channels, efficient administration, and precise accounting. The various administrative organs communicated with one another by means of a smoothly functioning network of messengers. Although almost 24,000 documents referring to the economy of Ur III have so far been published, the majority of them are still waiting to be properly evaluated. Nor is there yet a serviceable typology for them; only when that has been drawn up will it be possible to write a book entitled “The Economic System of Ur III.” Represented in the main by contracts (loans, leases of temple land, the purchase of slaves, and the like), the “private” sector makes up only a small part of this mass of textual material. Neither can the sites at which discoveries have been made so far be taken as representative. In northern Babylonia, for example, scarcely any contemporary written documents have yet been recovered.

Ethnic, geographic, and intellectual constituents. From the ethnic point of view, Mesopotamia was as heterogeneous at the end of the 3rd millennium as it had been earlier. The Akkadian element predominated, and the proportion of speakers of Akkadian to speakers of Sumerian continued to change in favour of the former. The third group, first mentioned under Shar-kali-sharri of Akkad, are the Amorites. In Ur III some members of this people are already found in the higher echelons of the administration, but most of them, organized in tribes, still led a nomadic life. Their great days came in the Old Babylonian period. While clearly differing linguistically from Akkadian, the Amorite language, which can be reconstructed to some extent from more than a thousand proper names, is fairly closely related to the so-called Canaanite branch of the Semitic languages, of which it may in fact represent an older form. The fact that King Shu-Sin had a regular wall built clear across the land, the “wall that keeps out the Tidnum” (the name of a tribe), shows how strong the pressure of the nomads was in the 21st century and what efforts were made to check their influx. The fourth major ethnic group was the Hurrians, who were especially important in northern Mesopotamia and in the vicinity of modern Kirkūk.

It is likely that the geographic horizon of the empire of Ur III did not materially exceed that of the empire of Akkad. No names of localities in the interior of Anatolia have been found, but there was much coming and going of messengers between Mesopotamia and Iran, far beyond Elam. There is also one mention of Gubla (Byblos) on the Mediterranean coast. Oddly enough, there is no evidence of any relations with Egypt, either in Ur III or in the Old Babylonian period. It is odd if no contacts existed

The Code
of Ur-
Nammu

The
Amorites

The *bala*
system



Sites associated with ancient Mesopotamian history.

at the end of the 3rd millennium between the two great civilizations of the ancient Middle East.

Intellectual life at the time of Ur III must have been very active in the cultivation and transmission of older literature, as well as in new creations. Although its importance as a spoken tongue was slowly diminishing, Sumerian still flourished as a written language, a state of affairs that continued into the Old Babylonian period. As shown by the hymn to the deified king, new literary genres arose in Ur III. If Old Babylonian copies are any indication, the king's correspondence with leading officials was also of a high literary level.

In the long view, the 3rd dynasty of Ur did not survive in historical memory as vigorously as did Akkad. To be sure, Old Babylonian historiography speaks of Ur III as *bala-Sulgi*, the "(reigning) cycle of Sulgi"; however, there is nothing that would correspond to the epic poems about Sargon and Naram-Sin. The reason is not clear, but it is conceivable that the later, purely Akkadian population felt a closer identification with Akkad than with a state that to a large extent still made use of the Sumerian language.

Ur III in decline. The decline of Ur III is an event in Mesopotamian history that can be followed in greater detail than other stages of that history thanks to sources such as the royal correspondence, two elegies on the destruction of Ur and Sumer, and an archive from Isin that shows how Ishbi-Erra, as usurper and king of Isin, eliminated his former overlord in Ur. Ibbi-Sin was waging war in Elam when an ambitious rival came forward in the person of Ishbi-Erra from Mari, presumably a general or high official. By emphasizing to the utmost the danger threatening from the Amorites, Ishbi-Erra urged the king to entrust to him the protection of the neighbouring cities of Isin and Nippur. Ishbi-Erra's demand came close to extortion, and his correspondence shows how skillfully he dealt with the Amorites and with individual *ensis*, some of whom soon went over to his side. Ishbi-Erra also took advantage of the depression that the king suffered because the god Enlil "hated him," a phrase presumably referring

to bad omens resulting from the examination of sacrificed animals, on which procedure many rulers based their actions (or, as the case may be, their inaction). Ishbi-Erra fortified Isin and, in the 10th year of Ibbi-Sin's reign, began to employ his own dating formulas on documents, an act tantamount to a renunciation of loyalty. Ishbi-Erra, for his part, believed himself to be the favourite of Enlil, the more so as he ruled over Nippur, where the god had his sanctuary. In the end he claimed suzerainty over all of southern Mesopotamia, including Ur.

While Ishbi-Erra purposefully strengthened his domains, Ibbi-Sin continued for 14 more years to rule over a decreasing portion of the land. The end of Ur came about through a concatenation of misfortunes: A famine broke out, and Ur was besieged, taken, and destroyed by the invading Elamites and their allies among the Iranian tribes. Ibbi-Sin was led away captive, and no more was heard of him. The elegies record in moving fashion the unhappy end of Ur, the catastrophe that had been brought about by the wrath of Enlil.

THE OLD BABYLONIAN PERIOD

Isin and Larsa. During the collapse of Ur III, Ishbi-Erra established himself in Isin and founded a dynasty there that lasted from 2017 to 1794. His example was followed elsewhere by local rulers, as in Dér, Eshnunna, Sippar, Kish, and Larsa. In many localities an urge was felt to imitate the model of Ur; Isin probably took over unchanged the administrative system of that state. Ishbi-Erra and his successors had themselves deified, as did one of the rulers of Dér, on the Iranian border. For almost a century Isin predominated within the mosaic of states that were slowly reemerging. Overseas trade revived after Ishbi-Erra had driven out the Elamite garrison from Ur, and under his successor, Shu-ilishu, a statue of the moon god Nanna, the city god of Ur, was recovered from the Elamites, who had carried it off. Up to the reign of Lipit-Ishtar (c. 1934–c. 1924), the rulers of Isin so resembled those of Ur, as far as the king's assessment of himself in the hymns is

Ishbi-Erra

concerned, that it seems almost arbitrary to postulate a break between Ibbi-Sin and Ishbi-Erra. As a further example of continuity it might be added that the Code of Lipit-Ishtar stands exactly midway chronologically between the Code of Ur-Nammu and the Code of Hammurabi. Yet it is much closer to the former in language and especially in legal philosophy than to Hammurabi's compilation of judgments. For example, the Code of Lipit-Ishtar does not know the *lex talionis* ("an eye for an eye and a tooth for a tooth"), the guiding principle of Hammurabi's penal law.

Political fragmentation. It is probable that the definitive separation from Ur III came about through changing components of the population, from "Sumerians and Akkadians" to "Akkadians and Amorites." An Old Babylonian liver omen states that "he of the steppes will enter, and chase out the one in the city." This is indeed an abbreviated formula for an event that took place more than once: the usurpation of the king's throne in the city by the "sheikh" of some Amorite tribe. These usurpations were regularly carried out as part of the respective tribes became settled, although this was not so in the case of Isin because the house of Ishbi-Erra came from Mari and was of Akkadian origin, to judge by the rulers' names. By the same linguistic token the dynasty of Larsa was Amorite. The fifth ruler of the latter dynasty, Gungunum (ruled c. 1932–c. 1906), conquered Ur and established himself as the equal and rival of Isin; at this stage—the end of the 20th century BC—if not before, Ur had certainly outlived itself. From Gungunum until the temporary unification of Mesopotamia under Hammurabi, the political picture was determined by the disintegration of the balance of power, by incessant vacillation of alliances, by the presumption of the various rulers, by the fear of encroachments by the Amorite nomads, and by increasingly wretched social conditions. The extensive archive of correspondence from the royal palace of Mari (c. 1810–1750) is the best source of information about the political and diplomatic game and its rules, whether honoured or broken; it covers treaties, the dispatch and reception of embassies, agreements about the integration of allied armies, espionage, and "situation reports" from "foreign" courts. Devoid of exaggeration or stylization, these letters, dealing as they do with everyday events, are preferable to the numerous royal inscriptions on buildings, even when the latter contain historical allusions.

Literary texts and increasing decentralization. Another indirect but far from negligible source for the political and socioeconomic situation in the 20th–18th centuries BC is the literature of omens. These are long compendiums in which the condition of a sheep's liver or some other divinatory object (for instance, the behaviour of a drop of oil in a beaker filled with water, the appearance of a newborn baby, and the shape of rising clouds of incense) is described at length and commented on with the appropriate prediction: "The king will kill his dignitaries and distribute their houses and property among the temples"; "A powerful man will ascend the throne in a foreign city"; "The land that rose up against its 'shepherd' will continue to be ruled by that 'shepherd'"; "The king will depose his chancellor"; and "They will lock the city gate and there will be a calamity in the city."

Beginning with Gungunum of Larsa, the texts allow greater insight into the private sector than in any other previous period. There is a considerable increase in the number of private contracts and private correspondences. Especially frequent among the private contracts are those concluded about loans of silver or grain (barley), illustrating the common man's plight, especially when driven to seek out a creditor, the first step on a road that in many instances led to ruin. The rate of interest, fixed at 20 percent in the case of silver and 33 percent in that of grain, increased further if the deadline for repayment, usually at harvest time, was not kept. Insolvency resulted in imprisonment for debt, slavery by mortgage, and even the sale of children and the debtor's own person. Many private letters contain entreaties for the release of family members from imprisonment at the creditor's hands. Yet considerable fortunes were also made, in "liquid" capital as well as landed property. As these tendencies threatened

to end in economic disaster, the kings prescribed as a corrective the liquidation of debts, by way of temporary alleviation at least. The exact wording of one such decree is known from the time of Ammišaduqa of Babylon.

Until the Ur III period, the only archives so far recovered dealt with temples or the palace. However, belonging to the Old Babylonian period, along with documents pertaining to civil law, were an increasing number of administrative records of privately managed households, inns, and farms: settlements of accounts, receipts, and notes on various transactions. Here was clearly a regular bourgeoisie, disposing of its own land and possessing means independent of temple and palace. Trade, too, was now chiefly in private hands; the merchant traveled (or sent his partners) at his own risk, not on behalf of the state. Among the civil-law contracts there was a substantial increase in records of land purchases. Also significant for the economic situation in the Old Babylonian era was a process that might be summarized as "secularization of the temples," even if all the stages of this development cannot be traced. The palace had probably possessed for centuries the authority to dispose of temple property, but, whereas Urukagina of Lagash had still branded the tendency as leading to abuses, the citizen's relationship to the temple now took on individual traits. Revenues from certain priestly offices—benefices, in other words—went to private individuals and were sold and inherited. The process had begun in Ur, where the king bestowed benefices, although the recipients could not own them. The archives of the "canonesses" of the sun god of Sippar furnish a particularly striking example of the fusion of religious service and private economic interest. These women, who lived in a convent called *gagûm*, came from the city's leading families and were not allowed to marry. With their property, consisting of land and silver, they engaged in a lively and remunerative business by granting loans and leasing out fields.

The tendency toward decentralization had begun in the Old Babylonian period with Isin. It concluded with the 72-year reign of the house of Kudur-Mabuk in Larsa (c. 1834–c. 1763). Kudur-Mabuk, sheikh of the Amorite tribe of the Jamutbal, despite his Elamite name, helped his son Warad-Sin to secure the throne. This usurpation allowed Larsa, which had passed through a period of internal unrest, to flourish one more time. Under Warad-Sin and in the long reign of his brother Rim-Sin, large portions of southern Babylonia, including Nippur, were once again united in one state of Larsa in 1794. Larsa was conquered by Hammurabi in 1763.

Early history of Assyria. Strictly speaking, the use of the name "Assyria" for the period before the latter half of the 2nd millennium BC is anachronistic; Assyria—as against the city-state of Ashur—did not become an independent state until about 1400 BC. For convenience, however, the term is used throughout this section.

In contrast to southern Mesopotamia or the mid-Euphrates region (Mari), written sources in Assyria do not begin until very late, shortly before Ur III. By Assyria—a region that does not lend itself to precise geographic delineation—is understood the territory on the Tigris north of the river's passage through the mountains of the Jabal Hamrîn to a point north of Nineveh, as well as the area between Little and Great Zab (a tributary of the Tigris in northeast Iraq) and to the north of the latter. In the north, Assyria was later bordered by the mountain state of Urartu; to the east and southeast its neighbour was the region around ancient Nuzi (near modern Kirkûk, "Arrapchitis" [Arrapkha] of the Greeks). In the early 2nd millennium the main cities of this region were Ashur (160 miles north-northwest of modern Baghdad), the capital (synonymous with the city god and national divinity); Nineveh, lying opposite modern Mosul; and Urbilum, later Arbela (modern Irbil, some 200 miles north of Baghdad).

In Assyria, inscriptions were composed in Akkadian from the beginning. Under Ur III, Ashur was a provincial capital. Assyria as a whole, however, is not likely to have been a permanently secured part of the empire, since two date formulas of Shulgi and Amar-Su'ena mention the destruction of Urbilum. Ideas of the population of Assyria in the 3rd millennium are necessarily very imprecise. It

Changing
relationship
between
religious
institutions
and
individuals

is not known how long Semitic tribes had been settled there. The inhabitants of southern Mesopotamia called Assyria Shubir in Sumerian and Subartu in Akkadian; these names may point to a Subarean population that was related to the Hurrians. Gasur, the later Nuzi, belonged to the Akkadian language region about the year 2200 but was lost to the Hurrians in the first quarter of the 2nd millennium. The Assyrian dialect of Akkadian found in the beginning of the 2nd millennium differs strongly from the dialect of Babylonia. These two versions of the Akkadian language continue into the 1st millennium.

In contrast to the kings of southern Mesopotamia, the rulers of Ashur styled themselves not king but partly *iš-šiakum*, the Akkadian equivalent of the Sumerian word *ensi*, partly *rubā'um*, or "great one." Unfortunately, the rulers cannot be synchronized precisely with the kings of southern Mesopotamia before Shamshi-Adad I (c. 1813–c. 1781 BC). For instance, it has not yet been established just when Ilushuma's excursion toward the southeast, recorded in an inscription, actually took place. Ilushuma boasts of having freed of taxes the "Akkadians and their children." While he mentions the cities of Nippur and Ur, the other localities listed were situated in the region east of the Tigris. The event itself may have taken place in the reign of Ishme-Dagan of Isin (c. 1953–c. 1935 BC), although how far Ilushuma's words correspond to the truth cannot be checked. In the Babylonian texts, at any rate, no reference is made to Assyrian intervention. The whole problem of dating is aggravated by the fact that the Assyrians did not, unlike the Babylonians, use date formulas that often contain interesting historical details; instead, every year was designated by the name of a high official (eponymic dating). The conscious cultivation of an old tradition is mirrored in the fact that two rulers of 19th-century Assyria called themselves Sargon and Naram-Sin after famous models in the Akkadian dynasty.

Aside from the generally scarce reports on projected construction, there is at present no information about the city of Ashur and its surroundings. There exists, however, unexpectedly rewarding source material from the trading colonies of Ashur in Anatolia. The texts come mainly from Kanesh (modern Kültepe, near Kayseri, in Turkey) and from Hattusa (modern Boğazköy, Tur.), the later Hittite capital. In the 19th century BC three generations of Assyrian merchants engaged in a lively commodity trade (especially in textiles and metal) between the homeland and Anatolia, also taking part profitably in internal Anatolian trade. Like their contemporaries in southern Mesopotamia, they did business privately and at their own risk, living peacefully and occasionally intermarrying with the "Anatolians." As long as they paid taxes to the local rulers, the Assyrians were given a free hand.

Clearly these forays by Assyrian merchants led to some transplanting of Mesopotamian culture into Anatolia. Thus the Anatolians adopted cuneiform writing and used the Assyrian language. While this influence doubtless already affected the first Hittites arriving in Anatolia, a direct line from the period of these trading colonies to the Hittite empire cannot yet be traced.

From about 1813 to about 1781 Assyria was ruled by Shamshi-Adad I, a contemporary of Hammurabi and a personality in no way inferior to him. Shamshi-Adad's father—an Amorite, to judge by the name—had ruled near Mari. The son, not being of Assyrian origin, ascended the throne of Assyria as a foreigner and on a detour, as it were, after having spent some time as an exile in Babylonia. He had his two sons rule as viceroys, in Ekallātum on the Tigris and in Mari, respectively, until the older of the two, Ishme-Dagan, succeeded his father on the throne. Through the archive of correspondence in the palace at Mari, scholars are particularly well informed about Shamshi-Adad's reign and many aspects of his personality. Shamshi-Adad's state had a common border for some time with the Babylonia of Hammurabi. Soon after Shamshi-Adad's death, Mari broke away, regaining its independence under an Amorite dynasty that had been living there for generations; in the end, Hammurabi conquered and destroyed Mari. After Ishme-Dagan's death, Assyrian history is lost sight of for more than 100 years.

The Old Babylonian empire. *Political fortunes.* Hammurabi (c. 1792–c. 1750 BC) is surely the most impressive and by now the best-known figure of the ancient Middle East of the first half of the 2nd millennium BC. He owes his posthumous reputation to the great stela into which the Code of Hammurabi was carved and indirectly also to the fact that his dynasty has made the name of Babylon famous for all time. In much the same way in which pre-Sargonic Kish exemplified the non-Sumerian area north of Sumer and Akkad lent its name to a country and a language, Babylonia became the symbol of the whole country that the Greeks called Babylonia. This term is used anachronistically by Assyriologists as a geographic concept in reference to the period before Hammurabi. Originally the city's name was probably Babilla, which was reinterpreted in popular etymology as Bāb-ili ("Gate of the God").

The 1st dynasty of Babylon rose from insignificant beginnings. The history of the erstwhile province of Ur is traceable from about 1894 onward, when the Amorite Sumuabum came to power there. What is known of these events fits altogether into the modest proportions of the period when Mesopotamia was a mosaic of small states. Hammurabi played skillfully on the instrument of coalitions and became more powerful than his predecessors had been. Nonetheless, it was only in the 30th year of his reign, after his conquest of Larsa, that he gave concrete expression to the idea of ruling all of southern Mesopotamia by "strengthening the foundations of Sumer and Akkad," in the words of that year's dating formula. In the prologue to the Code of Hammurabi the king lists the following cities as belonging to his dominions: Eridu near Ur, Ur, Lagash and Girsu, Zabalam, Larsa, Uruk, Adab, Isin, Nippur, Keshi, Dilbat, Borsippa, Babylon itself, Kish, Malgium, Mashkan-shapir, Kutha, Sippar, Eshnunna in the Diyālā region, Mari, Tuttul on the lower Balikh (a tributary of the Euphrates), and finally Ashur and Nineveh. This was on a scale reminiscent of Akkad or Ur III. Yet Ashur and Nineveh cannot have formed part of this empire for long because at the end of Hammurabi's reign mention is made again of wars against Subartu—that is, Assyria.

Under Hammurabi's son Samsuiluna (c. 1749–c. 1712 BC) the Babylonian empire greatly shrank in size. Following what had almost become a tradition, the south rose up in revolt. Larsa regained its autonomy for some time, and the walls of Ur, Uruk, and Larsa were leveled. Eshnunna, which evidently had also seceded, was vanquished about 1730. Later chronicles mention the existence of a state in the Sealand, with its own dynasty (by "Sealand" is understood the marshlands of southern Babylonia). Knowledge of this new dynasty is unfortunately very vague, only one of its kings being documented in contemporary texts. About 1741 Samsuiluna mentions the Kassites for the first time; about 1726 he constructed a stronghold, "Fort Samsuiluna," as a bulwark against them on the Diyālā near its confluence with the Tigris.

Like the Gutians before them, the Kassites were at first prevented from entering Babylonia and pushed into the mid-Euphrates region; there, in the kingdom of Khana (centred on Mari and Terqa, both below the junction with the Khābūr River), a king appears with the Kassite name of Kashtiliashu, who ruled toward the end of the Babylonian dynasty. From Khana the Kassites moved south in small groups, probably as harvest workers. After the Hittite invasion under Mursilis I, who is said to have dethroned the last king of Babylon, Samsudtana, in 1595, the Kassites assumed the royal power in Babylonia. So far, the contemporary sources do not mention this epoch, and the question remains unresolved as to how the Kassite rulers named in king lists mesh with the end of the 2nd millennium BC.

Babylonian law. The Code of Hammurabi is the most frequently cited cuneiform document in specialized literature. Its first scholarly publication in 1902 led to the development of a special branch of comparative jurisprudence, the study of cuneiform law. Following the division made by the first editor, Jean-Vincent Scheil, the Code of Hammurabi contains 280 judgments, or "paragraphs," on civil and criminal law, dealing in the main with cases

Problems with the Assyrian chronology

Hammurabi

The Code of Hammurabi

from everyday life in such a manner that it becomes obvious that the “lawgiver” or compiler had no intention of covering all possible contingencies. In broad outline, the themes treated in the Code of Hammurabi are libel; corrupt administration of justice; theft, receiving stolen goods, robbery, looting, and burglary; murder, manslaughter, and bodily injury; abduction; judicature of tax lessees; liability for negligent damage to fields and crop damage caused by grazing cattle; illegal felling of palm trees; legal problems of trade enterprises, in particular, the relationship between the merchant and his employee traveling overland, and embezzlement of merchandise; trust monies; the proportion of interest to loan money; the legal position of the female publican; slavery and ransom, slavery for debt, runaway slaves, the sale and manumission of slaves, and the contesting of slave status; the rent of persons, animals, and ships and their respective tariffs, offenses committed by hired labourers, and the vicious bull; family law: the price of a bride, dowry, the married woman’s property, wife and concubine, and the legal position of the respective issue, divorce, adoption, the wet nurse’s contract, and inheritance; and the legal position of certain priestesses.

A similar if much shorter compendium of judgments, probably antedating that of Hammurabi by a generation or two, has been discovered in Eshnunna.

Hammurabi, who called his own work *dināt mišarim*, or “verdicts of the just order,” states in the epilogue that it was intended as legal aid for persons in search of advice. Whether these judgments were meant to have binding force in the sense of modern statutes, however, is a matter of controversy. The Code of Hammurabi differs in many respects from the Code of Lipit-Ishtar, which was written in Sumerian. Its most striking feature lies in the extraordinary severity of its penalties and in the principle of the *lex talionis*. The same attitude is reflected in various Old Babylonian contracts in which defaulters are threatened with bodily punishment. It is often said, and perhaps rightly so, that this severity, which so contrasts with Sumerian judicial tradition, can be traced back to the Amorite influence.

There is yet another way in which the Code of Hammurabi has given rise to much discussion. Many of its “paragraphs” vary according to whether the case concerns an *awilum*, a *muškēnum*, or a *wardum*. A threefold division of the populace had been postulated on the basis of these distinctions. The *wardum* is the least problematic: he is the slave—that is, a person in bondage who could be bought and sold, unless he was able to regain his freedom under certain conditions as a debtor-slave. The *muškēnum* were, under King Hammurabi at least, persons employed by the palace who could be given land in usufruct without receiving it as property. *Awilum* were the citizens who owned land in their own right and depended neither on the palace nor on the temple. As the Soviet scholar Igor M. Diakonov has pointed out, the distinction cannot have been very sharply drawn, because the classes *awilum* and *muškēnum* are not mutually exclusive: a man in high palace office could fairly easily purchase land as private property, whereas the free citizen who got into debt as a result of a bad harvest or some other misfortune had one foot in the slave class. Still unanswered is the question as to which segment of the population could be conscripted to do public works, a term that included the levy in case of war.

Ammišaduqa (c. 1646–c. 1626 BC) comes a century and a half after Hammurabi. His edict, already referred to, lists, among others, the following social and economic factors: private debts in silver and grain, if arising out of loans, were canceled; also canceled were back taxes that certain officials owed the palace and that had to be collected from the people; the female publican had to renounce the collection of outstanding debts in beer and barley and was, in turn, excused from paying amounts of silver and barley to the king; taxes on leased property were reduced; debt slaves who had formerly been free (as against slaves made over from debtor to creditor) were ransomed; and high officials were forbidden on pain of death to press those who held property in fee into harvest work by prepayment of wages. The phrase “because the king gave the

land a just order” serves as a rationale for many of these instances. In contrast to the codes, about whose binding force there is much doubt, edicts such as those of Ammīšaduqa had legal validity since there are references to the edicts of other kings in numerous legal documents of the Old Babylonian period.

Babylonian literature. The literature and the literary languages of Babylonia during the three centuries following Ur III deserve attention. When commenting on literary and historical texts such as the inscriptions of the kings of Akkad, it was pointed out that these were not originals but copies of Old Babylonian vintage. So far, such copies are the main source for Sumerian literature. Yet, while the Old Babylonian period witnessed the creation of much literature (royal hymns of the kings of Isin, Larsa, and Babylon and elegies), it was above all a time of intensive cultivation of traditional literature. The great Sumerian poems, whose origins or first written version, respectively, can now be traced back to about 2600, were copied again and again. After 2000, when Sumerian as a spoken language rapidly receded to isolated regions and eventually disappeared altogether, texts began to be translated, line by line, into Akkadian until there came to be bilingual versions. An important part of this, especially in the instructional program in schools, were the so-called lexicographical texts. Sumerian word lists are almost as old as cuneiform writing itself; they formed the perfect material for those learning to write. In the Old Babylonian period, the individual lexical entries were translated and often annotated with phonetic signs. This led to the creation of “dictionaries,” the value of which to the modern philologist cannot be exaggerated. Since Sumerian had to be taught much more than before, regular “grammatical treatises” also came into being: so far as it was possible, in view of the radically different structures of the two languages, Sumerian pronouns, verb forms, and the like were translated into Akkadian, including entire “paradigms” of individual verbs.

In belles lettres, Sumerian still predominates, although there is no lack of Akkadian masterpieces, including the oldest Akkadian version of the epic of Gilgamesh. The very high prestige still enjoyed by Sumerian should not be underestimated, and it continued to be used for inscriptions on buildings and the yearly dating formulas. Aside from being the language of practical affairs (*i.e.*, letters and contracts), there was a high incidence of Akkadian in soothsaying and divinatory literature. To be sure, the Sumerians also practiced foretelling the future from the examination of animal entrails, but as far as is known they did not write down the results. In Akkadian, on the other hand, there are extensive and “scientifically” arranged compendiums of omens based on the liver (as well as other omens), reflecting the importance that the divination of the future had in religion, in politics, and in all aspects of daily life.

Judging by its increasingly refined juridical thought, its ability to master in writing ever more complicated administrative procedures, its advanced knowledge of mathematics, and the fact that it marks the beginning of the study of astronomy, the Old Babylonian period appears to have been a time of exceedingly active intellectual endeavour—despite, if not because of, its lack of political cohesiveness.

The Hurrians. The Hurrians enter the orbit of ancient Middle Eastern civilization toward the end of the 3rd millennium BC. They arrived in Mesopotamia from the north or the east, but it is not known how long they had lived in the peripheral regions. There is a brief inscription in Hurrian language from the end of the period of Akkad, while that of King Arishen (or Atalshen) of Urkish and Nawar is written in Akkadian. The language of the Hurrians must have belonged to a widespread group of ancient Middle Eastern languages. The relationship between Hurrian and Subarean has already been mentioned, and the language of the Uartians, who played an important role from the end of the 2nd millennium to the 8th century BC, is likewise closely related to Hurrian. According to the Soviet scholars Igor M. Diakonov and Sergei A. Starostin, the Eastern Caucasian languages are an offshoot of the Hurrian-Urartian group.

Literary
texts

Social
categories

It is not known whether the migrations of the Hurrians ever took the form of aggressive invasion; 18th-century-BC texts from Mari speak of battles with the Hurrian tribe of Turukku south of Lake Urmia (some 150 miles from the Caspian Sea's southwest corner), but these were mountain campaigns, not the warding off of an offensive. Proper names in cuneiform texts, their frequency increasing in the period of Ur III, constitute the chief evidence for the presence of Hurrians. Nevertheless, there is no clear indication that the Hurrians had already advanced west of the Tigris at that time. An entirely different picture results from the 18th-century palace archives of Mari and from texts originating near the upper Khābūr River. Northern Mesopotamia, west of the Tigris, and Syria appear settled by a population that is mainly Amorite and Hurrian; and the latter had already reached the Mediterranean littoral, as shown by texts from Alalakh on the Orontes. In Mari, literary texts in Hurrian also have been found, indicating that Hurrian had by then become a fully developed written language as well.

The high point of the Hurrian period was not reached until about the middle of the 2nd millennium. In the 15th century, Alalakh was heavily Hurrianized; and in the empire of Mitanni the Hurrians represented the leading and perhaps the most numerous population group. (D.O.E.)

Mesopotamia to the end of the Achaemenian period

THE KASSITES, THE MITANNI, AND THE RISE OF ASSYRIA

About 150 years after the death of Hammurabi, his dynasty was destroyed by an invasion of new peoples. Because there are very few written records from this era, the time from about 1560 BC to about 1440 BC (in some areas until 1400 BC) is called the dark ages. The remaining Semitic states, such as the state of Ashur, became minor states within the sphere of influence of the new states of the Kassites and the Hurrians/Mitanni. The languages of the older cultures, Akkadian and Sumerian, continued or were soon reestablished, however. The cuneiform script persisted as the only type of writing in the entire area. Cultural continuity was not broken off, either, particularly in Babylonia. A matter of importance was the emergence of new Semitic leading classes from the ranks of the priesthood and the scribes. These gained increasing power.

The Kassites in Babylonia. The Kassites had settled by 1800 BC in what is now western Iran in the region of Hamadan-Kermanshah. The first to feel their forward thrust was Samsuiluna, who had to repel groups of Kassite invaders. Increasing numbers of Kassites gradually reached Babylonia and other parts of Mesopotamia. There they founded principalities, of which little is known. No inscription or document in the Kassite language has been preserved. Some 300 Kassite words have been found in Babylonian documents. Nor is much known about the social structure of the Kassites or their culture. There seems to have been no hereditary kingdom. Their religion was polytheistic; the names of some 30 gods are known.

The beginning of Kassite rule in Babylonia cannot be dated exactly. A king called Agum II ruled over a state that stretched from western Iran to the middle part of the Euphrates valley; 24 years after the Hittites had carried off the statue of the Babylonian god Marduk, he regained possession of the statue, brought it back to Babylon, and renewed the cult, making the god Marduk the equal of the corresponding Kassite god, Shuqamuna. Meanwhile, native princes continued to reign in southern Babylonia. It may have been Ulamburiash who finally annexed this area around 1450 and began negotiations with Egypt in Syria. Karaindash built a temple with bas-relief tile ornaments in Uruk (Erech) around 1420. A new capital west of Baghdad, Dūr Kurigalzu, competing with Babylon, was founded and named after Kurigalzu I (c. 1400–c. 1375). His successors Kadashman-Enlil I (c. 1375–c. 1360) and Burnaburiash II (c. 1360–c. 1333) were in correspondence with the Egyptian rulers Amenhotep III and Akhenaton (Amenhotep IV). They were interested in trading their lapis lazuli and other items for gold as well as in planning political marriages. Kurigalzu II (c. 1332–c. 1308) fought

against the Assyrians but was defeated by them. His successors sought to ally themselves with the Hittites in order to stop the expansion of the Assyrians. During the reign of Kashtiliash IV (c. 1232–c. 1225), Babylonia waged war on two fronts at the same time—against Elam and Assyria—ending in the catastrophic invasion and destruction of Babylon by Tukulti-Ninurta I. Not until the time of the kings Adad-shum-ušur (c. 1216–c. 1187) and Meli-shipak (c. 1186–c. 1172) was Babylon able to experience a period of prosperity and peace. Their successors were again forced to fight, facing the conqueror King Shutruk-Nahhunte of Elam (c. 1185–c. 1155). Cruel and fierce, the Elamites finally destroyed the dynasty of the Kassites during these wars (about 1155). Some poetical works lament this catastrophe.

Letters and documents of the time after 1380 show that many things had changed after the Kassites took power. The Kassite upper class, always a small minority, had become largely "Babylonianized." Babylonian names were to be found even among the royalty, and they predominated among the civil servants and the officers. The new feudal character of the social structure showed the influence of the Kassites. Babylonian town life had revived on the basis of commerce and handicrafts. The Kassitic nobility, however, maintained the upper hand in the rural areas, their wealthiest representatives holding very large landed estates. Many of these holdings came from donations of the king to deserving officers and civil servants, considerable privileges being connected with such grants. From the time of Kurigalzu II these were registered on stone tablets or, more frequently, on boundary stones called *kudurrus*. After 1200 the number of these increased substantially, because the kings needed a steadily growing retinue of loyal followers. The boundary stones had pictures in bas-relief, very often a multitude of religious symbols, and frequently contained detailed inscriptions giving the borders of the particular estate; sometimes the deserts of the recipient were listed and his privileges recorded; finally, trespassers were threatened with the most terrifying curses. Agriculture and cattle husbandry were the main pursuits on these estates, and horses were raised for the light war chariots of the cavalry. There was an export trade in horses and vehicles in exchange for raw material. As for the king, the idea of the social-minded ruler continued to be valid.

The decline of Babylonian culture at the end of the Old Babylonian period continued for some time under the Kassites. Not until approximately 1420 did the Kassites develop a distinctive style in architecture and sculpture. Kurigalzu I played an important part, especially in Ur, as a patron of the building arts. Poetry and scientific literature developed only gradually after 1400. The existence of earlier work is clear from poetry, philological lists, and collections of omens and signs that were in existence by the 14th century or before and that have been discovered in the Hittite capital of Hattusa, in the Syrian capital of Ugarit, and even as far away as Palestine. Somewhat later, new writings appear: medical diagnoses and recipes, more Sumerian-Akkadian word lists, and collections of astrological and other omens and signs with their interpretations. Most of these works are known today only from copies of more recent date. The most important is the Babylonian epic of the creation of the world, *Enuma elish*. Composed by an unknown poet, probably in the 14th century, it tells the story of the god Marduk. He began as the god of Babylon and was elevated to be king over all other gods after having successfully accomplished the destruction of the powers of chaos. For almost 1,000 years this epic was recited during the New Year's festival in the spring as part of the Marduk cult in Babylon. The literature of this time contains very few Kassitic words. Many scholars believe that the essential groundwork for the development of the subsequent Babylonian culture was laid during the later epoch of the Kassite era.

The Hurrian and Mitanni kingdoms. The weakening of the Semitic states in Mesopotamia after 1550 enabled the Hurrians to penetrate deeper into this region, where they founded numerous small states in the eastern parts of Anatolia, Mesopotamia, and Syria. The Hurrians came from northwestern Iran, but until recently very little was known

Society
under the
Kassites

The dark
ages

Possible
Aryan
influences

about their early history. After 1500, isolated dynasties appeared with Indo-Aryan names, but the significance of this is disputed. The presence of Old Indian technical terms in later records about horse breeding and the use of the names of Indian gods (such as, for example, Indra and Varuna) in some compacts of state formerly led several scholars to assume that numerous groups of Aryans, closely related to the Indians, pushed into Anatolia from the northeast. They were also credited with the introduction of the light war chariot with spoked wheels. This conclusion, however, is by no means established fact. So far it has not been possible to appraise the numbers and the political and cultural influence of the Aryans in Anatolia and Mesopotamia relative to those of the Hurrians.

Some time after 1500 the kingdom of Mitanni (or Mittani) arose near the sources of the Khābūr River in Mesopotamia. Since no record or inscription of their kings has been unearthed, little is known about the development and history of the Mitanni kingdom before King Tushratta. The Mitanni empire was known to the Egyptians under the name of Naharina, and Thutmose III fought frequently against it after 1460 BC. By 1420 the domain of the Mitanni king Saustatar (Saushatar) stretched from the Mediterranean all the way to the northern Zagros Mountains, in western Iran, including Alalakh, in northern Syria, as well as Nuzi, Kurrukhanni, and Arrapkha. The northern boundary dividing Mitanni from the Hittites and the other Hurrian states was never fixed, even under Saustatar's successors Artatama I and Shuttarna II, who married their daughters to the pharaohs Thutmose IV (1400–1390) and Amenhotep III (1390–1353). Tushratta (c. 1365–c. 1330), the son of Shuttarna, was able to maintain the kingdom he had inherited for many years. In his sometimes very long letters—one of them written in Hurrian—to Amenhotep III and Akhenaton (1353–1336), he wrote about commerce, his desire for gold, and marriage. Weakened by internal strife, the Mitanni kingdom eventually became a pawn between the rising kingdoms of the Hittites and the Assyrians.

The kingdom of Mitanni was a feudal state led by a warrior nobility of Aryan or Hurrian origin. Frequently horses were bred on their large landed estates. Documents and contract agreements in Syria often mention a chariot-warrior caste that also constituted the social upper class in the cities. The aristocratic families usually received their landed property as an inalienable fief. Consequently, no documents on the selling of landed property are to be found in the great archives of Akkadian documents and letters discovered in Nuzi, near Kirkūk. The prohibition against selling landed property was often dodged, however, with a stratagem: the previous owner “adopted” a willing buyer against an appropriate sum of money. The wealthy lord Tehiptilla was “adopted” almost 200 times, acquiring tremendous holdings of landed property in this way without interference by the local governmental authorities. He had gained his wealth through trade and commerce and through a productive two-field system of agriculture (in which each field was cultivated only once in two years). For a long time, Prince Shilwa-Teshub was in charge of the royal governmental administration in the district capital. Sheep breeding was the basis for a woolen industry, and textiles collected by the palace were exported on a large scale. Society was highly structured in classes, ranks, and professions. The judiciary, patterned after the Babylonian model, was well organized; the documents place heavy emphasis on correct procedure.

Native sources on the religion of the Hurrians of the Mitanni kingdom are limited; about their mythology, however, much is known from related Hittite and Ugaritic myths. Like the other peoples of the ancient Middle East, the Hurrians worshiped gods of various origins. The king of the gods was the weather god Teshub. According to the myths, he violently deposed his father Kumarbi; in this respect he resembled the Greek god Zeus, who deposed his father Kronos. The war chariot of Teshub was drawn by the bull gods Seris (“Day”) and Hurris (“Night”). Major sanctuaries of Teshub were located at Arrapkha (modern Kirkūk) and at Halab (modern Aleppo) in Syria. In the east his consort was the goddess of love and war Shaushka,

Hurrian
religion

and in the west the goddess Hebat (Hepat); both were similar to the Ishtar-Astarte of the Semites.

The sun god Shimegi and the moon god Kushuh, whose consort was Nikkal, the Ningal of the Sumerians, were of lesser rank. More important was the position of the Babylonian god of war and the underworld, Nergal. In northern Syria the god of war Astapi and the goddess of oaths Ishara are attested as early as the 3rd millennium BC.

In addition, a considerable importance was attributed to impersonal numina such as heaven and earth as well as to deities of mountains and rivers. In the myths the terrible aspect of the gods often prevails over indications of a benevolent attitude. The cults of sacrifices and other rites are similar to those known from the neighbouring countries; many Hurrian rituals were found in Hittite Anatolia. There is abundant evidence for magic and oracles.

Temple monuments of modest dimensions have been unearthed; in all probability, specific local traditions were a factor in their design. The dead were probably buried outside the settlement. Small artifacts, particularly seals, show a peculiar continuation of Babylonian and Assyrian traditions in their preference for the naturalistic representation of figures. There were painted ceramics with finely drawn decorations (white on a dark background). The strong position of the royal house was evident in the large palaces, existing even in district capitals. The palaces were decorated with frescoes. Because only a few Mitanni settlements have been unearthed in Mesopotamia, knowledge of Mitanni arts and culture is as yet insufficient.

The rise of Assyria. Very little can be said about northern Assyria during the 2nd millennium BC. Information on the old capital, Ashur, located in the south of the country, is somewhat more plentiful. The old lists of kings suggest that the same dynasty ruled continuously over Ashur from about 1600. All the names of the kings are given, but little else is known about Ashur before 1420. Almost all the princes had Akkadian names, and it can be assumed that their sphere of influence was rather small. Although Assyria belonged to the kingdom of the Mitanni for a long time, it seems that Ashur retained a certain autonomy. Located close to the boundary with Babylonia, it played that empire off against Mitanni whenever possible. Puzur-Ashur III concluded a border treaty with Babylonia about 1480, as did Ashur-bel-nisheshu about 1405. Ashur-nadin-ahhe II (c. 1392–c. 1383) was even able to obtain support from Egypt, which sent him a consignment of gold.

Ashur-uballit I (c. 1354–c. 1318) was at first subject to King Tushratta of Mitanni. After 1340, however, he attacked Tushratta, presumably together with Suppiluliumas I of the Hittites. Taking away from Mitanni parts of northeastern Mesopotamia, Ashur-uballit now called himself “Great King” and socialized with the king of Egypt on equal terms, arousing the indignation of the king of Babylonia. Ashur-uballit was the first to name Assyria the Land of Ashur, because the old name, Subartu, was often used in a derogatory sense in Babylonia. He ordered his short inscriptions to be partly written in the Babylonian dialect rather than the Assyrian, since this was considered refined. Marrying his daughter to a Babylonian, he intervened there energetically when Kassite nobles murdered his grandson. Future generations came to consider him rightfully as the real founder of the Assyrian empire. His son Enlil-nirari (c. 1326–c. 1318) also fought against Babylonia. Arik-den-ili (c. 1308–c. 1297) turned westward, where he encountered Semitic tribes of the so-called Akhlamu group.

Still greater successes were achieved by Adad-nirari I (c. 1295–c. 1264). Defeating the Kassite king Nazimaruttash, he forced him to retreat. After that he defeated the kings of Mitanni, first Shattuara I, then Wasashatta. This enabled him for a time to incorporate all Mesopotamia into his empire as a province, although in later struggles he lost large parts to the Hittites. In the east, he was satisfied with the defense of his lands against the mountain tribes.

Adad-nirari's inscriptions were more elaborate than those of his predecessors and were written in the Babylonian dialect. In them he declares that he feels called to these wars by the gods, a statement that was to be repeated by other kings after him. Assuming the old title of great king, he

Assyria
under
Mitanni
rule

called himself "King of All." He enlarged the temple and the palace in Ashur and also developed the fortifications there, particularly at the banks of the Tigris River. He worked on large building projects in the provinces.

Shalmaneser I

His son Shalmaneser I (Shulmanu-asharidu; c. 1263–c. 1234) attacked Uruartu (later called Urartu) in southern Armenia, which had allegedly broken away. Shattuara II of Hanigalbat, however, put him into a difficult situation, cutting his forces off from their water supplies. With courage born of despair, the Assyrians fought themselves free. They then set about reducing what was left of the Mitanni kingdom into an Assyrian province. The king claimed to have blinded 14,400 enemies in one eye—psychological warfare of a similar kind was used more and more as time went by. The Hittites tried in vain to save Hanigalbat. Together with the Babylonians they fought a commercial war against Ashur for many years. Like his father, Shalmaneser was a great builder. At the juncture of the Tigris and Great Zab rivers, he founded a strategically situated second capital, Kalakh (biblical Calah; modern Nimrūd).

His son was Tukulti-Ninurta (c. 1233–c. 1197), the Ninus of Greek legends. Gifted but extravagant, he made his nation a great power. He carried off thousands of Hittites from eastern Anatolia. He fought particularly hard against Babylonia, deporting Kashtiliash IV to Assyria. When the Babylonians rebelled again, he plundered the temples in Babylon, an act regarded as a sacrilege, even in Assyria. The relationship between the king and his capital deteriorated steadily. For this reason the king began to build a new city, Kar-Tukulti-Ninurta, on the other side of the Tigris River. Ultimately, even his sons rebelled against him and laid siege to him in his city; in the end he was murdered. His victorious wars against Babylonia were glorified in an epic poem, but his empire broke up soon after his death. Assyrian power declined for a time, while that of Babylonia rose.

Assyria had suffered under the oppression of both the Hurrians and the Mitanni kingdom. Its struggle for liberation and the bitter wars that followed had much to do with its development into a military power. In his capital of Ashur, the king depended on the citizen class and the priesthood, as well as on the landed nobility that furnished him with the war-chariot troops.

Documents and letters show the important role that agriculture played in the development of the state. Assyria was less dependent on artificial irrigation than was Babylonia. The breeding of horses was carried on intensively; remnants of elaborate directions for their training are extant. Trade and commerce also were of notable significance: metals were imported from Anatolia or Armenia, tin from northwestern Iran, and lumber from the west. The opening up of new trade routes was often a cause and the purpose of war.

Assyrian architecture, derived from a combination of Mitannian and Babylonian influences, developed early quite an individual style. The palaces often had colourful wall decorations. The art of seal cutting, taken largely from Mitanni, continued creatively on its own. The schools for scribes, where all the civil servants were trained, taught both the Babylonian and the Assyrian dialects of the Akkadian language. Babylonian works of literature were assimilated into Assyrian, often reworked into a different form. The Hurrian tradition remained strong in the military and political sphere while at the same time influencing the vocabulary of language.

Artistic and cultural developments

ASSYRIA AND BABYLONIA AT THE END OF THE 2ND MILLENNIUM

Babylonia under the 2nd dynasty of Isin. In a series of heavy wars about which not much is known, Marduk-kabit-ahheshu (c. 1152–c. 1135) established what came to be known as the 2nd dynasty of Isin. His successors were often forced to continue the fighting. The most famous king of the dynasty was Nebuchadnezzar I (Nabu-kudurriuşur; c. 1119–c. 1098). He fought mainly against Elam, which had conquered and ravaged a large part of Babylonia. His first attack miscarried because of an epidemic among his troops, but in a later campaign he conquered

Susa, the capital of Elam, and returned the previously removed statue of the god Marduk to its proper place. Soon thereafter the king of Elam was assassinated, and his kingdom once again fell apart into small states. This enabled Nebuchadnezzar to turn west, using the later years of peace to start extensive building projects. After him, his son became king, succeeded by his brother Marduk-nadin-ahhe (c. 1093–c. 1076). At first successful in his wars against Assyria, he later experienced heavy defeat. A famine of catastrophic proportions triggered an attack from Aramaean tribes, the ultimate blow. His successors made peace with Assyria, but the country suffered more and more from repeated attacks by Aramaeans and other Semitic nomads. Even though some of the kings still assumed grand titles, they were unable to stem the progressive disintegration of their empire. There followed the era known as the 2nd dynasty of the Sealand (c. 1020–c. 1000), which included three usurpers. The first of these had the Kassitic name of Simbar-Shihu (or Simbar-Shipak; c. 1020–c. 1003).

Toward the end of its reign, the dynasty of the Kassites became completely Babylonianized. The changeover to the dynasty of Isin, actually a succession of kings from different families, brought no essential transformation of the social structure. The feudal order remained. New landed estates came into existence in many places through grants to deserving officers; many boundary stones (*kudurrus*) have been found that describe them. The cities of Babylonia retained much of their former autonomy. The border provinces, however, were administered by royally appointed governors with civil and military functions.

In the literary arts this was a period of creativity; thus the later Babylonians with good reason regarded the time of Nebuchadnezzar I as one of the great eras of their history. A heroic epic, modeled upon older epics, celebrates the deeds of Nebuchadnezzar I, but unfortunately little of it is extant. Other material comes from the ancient myths. The poet of the later version of the epic of Gilgamesh, Sin-leqe-unnini (c. 1150–?) of Uruk, is known by name. This version of the epic is known as the Twelve-Tablet Poem; it contains about 3,000 verses. It is distinguished by its greater emphasis on the human qualities of Gilgamesh and his friend Enkidu; this quality makes it one of the great works of world literature.

The epic of Gilgamesh

Another poet active at about the same time was the author of a poem of 480 verses called *Ludlul bēl nēmeqi* ("Let Me Praise the Possessor of Wisdom"). The poem meditates on the workings of divine justice, which sometimes appear strange and inexplicable to suffering human beings; this subject had acquired an increasing importance in the contemporary religion of Babylon. The poem describes the multifarious sufferings of a high official and his subsequent salvation by the god Marduk.

The gradual reduction of the Sumerian pantheon of about 2,000 gods by the identification and integration of originally distinct gods and goddesses of similar functions resulted in a growing number of surnames or compound names for the main gods (Marduk, for example, had about 50 such names) and later in a conception of "the god" and "the goddess" with interchangeable names in the cults of the great temples. There was a theology of identifications of gods, which was documented by god lists in two columns with hundreds of entries in the form "En-zag = Nabū of (the island of) Dilmun," as well as by many hymns and prayers of the time and by later compositions.

As a consequence of the distinction of an enormous number of multifarious sins, the concept of a universal sinfulness of mankind is increasingly observed in this period and later. All human beings, therefore, were believed to be in need of the forgiveness afforded by the deities to sincere worshipers. Outside of Israel, the concept of sinfulness can be found in ancient times only in Babylonia and Assyria.

Assyria between 1200 and 1000 BC. After a period of decline following Tukulti-Ninurta I, Assyria was consolidated and stabilized under Ashur-dan I (c. 1179–c. 1134) and Ashur-resh-ishi I (c. 1133–c. 1116). Several times forced to fight against Babylonia, the latter was even able to defend himself against an attack by Nebuchadnezzar I.

According to the inscriptions, most of his building efforts were in Nineveh, rather than in the old capital of Ashur.

His son Tiglath-pileser I (Tukulti-apil-Esharra; c. 1115–c. 1077) raised the power of Assyria to new heights. First he turned against a large army of the Mushki that had entered into southern Armenia from Anatolia, defeating them decisively. After this, he forced the small Hurrian states of southern Armenia to pay him tribute. Trained in mountain warfare themselves and helped by capable pioneers, the Assyrians were now able to advance far into the mountain regions. Their main enemies were the Aramaeans, the Semitic Bedouin nomads whose many small states often combined against the Assyrians. Tiglath-pileser I also went to Syria and even reached the Mediterranean, where he took a sea voyage. After 1100 these campaigns led to conflicts with Babylonia. Tiglath-pileser conquered northern Babylonia and plundered Babylon, without decisively defeating Marduk-nadin-ahhe. In his own country the king paid particular attention to agriculture and fruit growing, improved the administrative system, and developed more thorough methods of training scribes.

Three of his sons reigned after Tiglath-pileser, including Ashur-bel-kala (c. 1074–c. 1057). Like his father, he fought in southern Armenia and against the Aramaeans with Babylonia as his ally. Disintegration of the empire could not be delayed, however. The grandson of Tiglath-pileser, Ashurnasirpal I (c. 1050–c. 1032), was sickly and unable to do more than defend Assyria proper against his enemies. Fragments of three of his prayers to Ishtar are preserved; among them is a penitential prayer in which he wonders about the cause of so much adversity. Referring to his many good deeds but admitting his guilt at the same time, he asks for forgiveness and health. According to the king, part of his guilt lay in neglecting to teach his subjects the fear of god. After him, little is known for 100 years.

State and society during the time of Tiglath-pileser were not essentially different from those of the 13th century. Collections of laws, drafts, and edicts of the court exist that go back as far as the 14th century BC. Presumably, most of these remained in effect. One tablet defining the marriage laws shows that the social position of women in Assyria was lower than in Babylonia or Israel or among the Hittites. A man was allowed to send away his wife at his own pleasure with or without divorce money. In the case of adultery, he was permitted to kill or maim her. Outside her house the woman was forced to observe many restrictions, such as the wearing of a veil. It is not clear whether these regulations carried the weight of law, but they seem to have represented a reaction against practices that were more favourable to women. Two somewhat older marriage contracts, for example, granted equal rights to both partners, even in divorce. The women of the king's harem were subject to severe punishment, including beating, maiming, and death, along with those who guarded and looked after them. The penal laws of the time were generally more severe in Assyria than in other countries of the East. The death penalty was not uncommon. In less serious cases the penalty was forced labour after flogging. In certain cases there was trial by ordeal. One tablet treats the subject of landed property rights. Offenses against the established boundary lines called for extremely severe punishment. A creditor was allowed to force his debtor to work for him, but he could not sell him.

The greater part of Assyrian literature was either taken over from Babylonia or written by the Assyrians in the Babylonian dialect, who modeled their works on Babylonian originals. The Assyrian dialect was used in legal documents, court and temple rituals, and collections of recipes—as, for example, in directions for making perfumes. A new art form was the picture tale: a continuing series of pictures carved on square stelae of stone. The pictures, showing war or hunting scenes, begin at the top of the stela and run down around it, with inscriptions under the pictures explaining them. These and the finely cut seals show that the fine arts of Assyria were beginning to surpass those of Babylonia. Architecture and other forms of the monumental arts also began a further development, such as the double temple with its two towers (ziggurat). Colourful enameled tiles were used to decorate the facades.

ASSYRIA AND BABYLONIA FROM C. 1000 TO C. 750 BC

Assyria and Babylonia until Ashurnasirpal II. The most important factor in the history of Mesopotamia in the 10th century was the continuing threat from the Aramaean seminomads. Again and again, the kings of both Babylonia and Assyria were forced to repel their invasions. Even though the Aramaeans were not able to gain a foothold in the main cities, there are evidences of them in many rural areas. Ashur-dan II (934–912) succeeded in suppressing the Aramaeans and the mountain people, in this way stabilizing the Assyrian boundaries. He reintroduced the use of the Assyrian dialect in his written records.

Adad-nirari II (c. 911–891) left detailed accounts of his wars and his efforts to improve agriculture. He led six campaigns against Aramaean intruders from northern Arabia. In two campaigns against Babylonia he forced Shamash-mudammiq (c. 930–904) to surrender extensive territories. Shamash-mudammiq was murdered, and a treaty with his successor, Nabu-shum-ukin (c. 904–888), secured peace for many years. Tukulti-Ninurta II (c. 890–884), the son of Adad-nirari II, preferred Nineveh to Ashur. He fought campaigns in southern Armenia. He was portrayed on stelae in blue and yellow enamel in the late Hittite style, showing him under a winged sun—a theme adopted from Egyptian art. His son Ashurnasirpal II (883–859) continued the policy of conquest and expansion. He left a detailed account of his campaigns, which were impressive in their cruelty. Defeated enemies were impaled, flayed, or beheaded in great numbers. Mass deportations, however, were found to serve the interests of the growing empire better than terror. Through the systematic exchange of native populations, conquered regions were denationalized. The result was a submissive, mixed population in which the Aramaean element became the majority. This provided the labour force for the various public works in the metropolitan centres of the Assyrian empire. Ashurnasirpal II rebuilt Kalakh, founded by Shalmaneser I, and made it his capital. Ashur remained the centre of the worship of the god Ashur—in whose name all the wars of conquest were fought. A third capital was Nineveh.

Ashurnasirpal II was the first to use cavalry units to any large extent in addition to infantry and war-chariot troops. He also was the first to employ heavy, mobile battering rams and wall breakers in his sieges. Following after the conquering troops came officials from all branches of the civil service, because the king wanted to lose no time in incorporating the new lands into his empire. The supremacy of Assyria over its neighbouring states owed much to the proficiency of the government service under the leadership of the minister Gabbilani-eresh. The campaigns of Ashurnasirpal II led him mainly to southern Armenia and Mesopotamia. After a series of heavy wars, he incorporated Mesopotamia as far as the Euphrates River. A campaign to Syria encountered little resistance. There was no great war against Babylonia. Ashurnasirpal, like other Assyrian kings, may have been moved by religion not to destroy Babylonia, which had almost the same gods as Assyria. Both empires must have profited from mutual trade and cultural exchange. The Babylonians, under the energetic Nabu-apla-iddina (c. 887–855) attacked the Aramaeans in southern Mesopotamia and occupied the valley of the Euphrates River to about the mouth of the Khäbür River.

Ashurnasirpal, so brutal in his wars, was able to inspire architects, structural engineers, and artists and sculptors to heights never before achieved. He built and enlarged temples and palaces in several cities. His most impressive monument was his own palace in Kalakh, covering a space of 269,000 square feet (25,000 square metres). Hundreds of large limestone slabs were used in murals in the state-rooms and living quarters. Most of the scenes were done in relief, but painted murals also have been found. Most of them depict mythological themes and symbolic fertility rites, with the king participating. Brutal war pictures were aimed to discourage enemies. The chief god of Kalakh was Ninurta, god of war and the hunt. The tower of the temple dedicated to Ninurta also served as an astronomical observatory. Kalakh soon became the cultural centre of the empire. Ashurnasirpal claimed to have entertained 69,574 guests at the opening ceremonies of his palace.

Aramaean
invasions

Disintegra-
tion of the
empire

Develop-
ments in
art and
architec-
ture under
Ashurnasir-
pal II

Shalmaneser III and Shamshi-Adad V of Assyria. The son and successor of Ashurnasirpal was Shalmaneser III (858–824). His father's equal in both brutality and energy, he was less realistic in his undertakings. His inscriptions, in a peculiar blend of Assyrian and Babylonian, record his considerable achievements but are not always able to conceal his failures. His campaigns were directed mostly against Syria. While he was able to conquer northern Syria and make it a province, in the south he could only weaken the strong state of Damascus and was unable, even after several wars, to eliminate it. In 841 he laid unsuccessful siege to Damascus. Also in 841 King Jehu of Israel was forced to pay tribute. In his invasion of Cilicia, Shalmaneser had only partial success. The same was true of the kingdom of Urartu in Armenia, from which, however, the troops returned with immense quantities of lumber and building stone. The king and, in later years, the general Dayyan-Ashur went several times to western Iran, where they found such states as Mannai in northwestern Iran and, farther away in the southeast, the Persians. They also encountered the Medes during these wars. Horse tribute was collected.

In Babylonia, Marduk-zakir-shumi I ascended the throne about the year 855. His brother Marduk-bel-usati rebelled against him, and in 851 the king was forced to ask Shalmaneser for help. Shalmaneser was only too happy to oblige; when the usurper had been finally eliminated (850), Shalmaneser went to southern Babylonia, which at that time was almost completely dominated by Aramaeans. There he encountered, among others, the Chaldeans, mentioned for the first time in 878 BC, who were to play a leading role in the history of later times; Shalmaneser made them tributaries.

During his long reign he built temples, palaces, and fortifications in Assyria as well as in the other capitals of his provinces. His artists created many statues and stelae. Among the best known is the Black Obelisk, which includes a picture of Jehu of Israel paying tribute. The bronze doors from the town of Imgur-Enlil (Balawat) in Assyria portray the course of his campaigns and other undertakings in rows of pictures, often very lifelike. Hundreds of delicately carved ivories were carried away from Phoenicia, and many of the artists along with them; these later made Kalakh a centre for the art of ivory sculpture.

In the last four years of the reign of Shalmaneser, the crown prince Ashur-da'in-apla led a rebellion. The old king appointed his younger son Shamshi-Adad as the new crown prince. Forced to flee to Babylonia, Shamshi-Adad V (823–811) finally managed to regain the kingship with the help of Marduk-zakir-shumi I under humiliating conditions. As king he campaigned with varying success in southern Armenia and Azerbaijan, later turning against Babylonia. He won several battles against the Babylonian kings Marduk-balassu-iqbi and Baba-aha-iddina (about 818–12) and pushed through to Chaldea. Babylonia remained independent, however.

Adad-nirari III and his successors. Shamshi-Adad V died while Adad-nirari III (810–783) was still a minor. His Babylonian mother, Sammu-ramat, took over the regency, governing with great energy until 806. The Greeks, who called her Semiramis, credited her with legendary accomplishments, but historically little is known about her. Adad-nirari later led several campaigns against the Medes and also against Syria and Palestine. In 804 he reached Gaza, but Damascus proved invincible. He also fought in Babylonia, helping to restore order in the north.

Shalmaneser IV (c. 783–773) fought against Urartu, then at the height of its power under King Argishti (c. 780–755). He successfully defended eastern Mesopotamia against attacks from Armenia. On the other hand, he lost most of Syria after a campaign against Damascus in 773. The reign of Ashur-dan III (772–755) was shadowed by rebellions and by epidemics of plague. Of Ashur-nirari V (754–746) little is known.

In Assyria the feudal structure of society remained largely unchanged. Many of the conquered lands were combined to form large provinces. The governors of these provinces sometimes acquired considerable independence, particularly under the weaker monarchs after Adad-nirari III.

Some of them even composed their own inscriptions. The influx of displaced peoples into the cities of Assyria created large metropolitan centres. The spoils of war, together with an expanding trade, favoured the development of a well-to-do commercial class. The dense population of the cities gave rise to social tensions that only the strong kings were able to contain. A number of the former capitals of the conquered lands remained important as capitals of provinces. There was much new building. A standing occupational force was needed in the provinces, and these troops grew steadily in proportion to the total military forces. There are no records on the training of officers or on military logistics. The civil service also expanded, the largest administrative body being the royal court, with thousands of functionaries and craftsmen in the several residential cities.

The cultural decline about the year 1000 was overcome during the reigns of Ashurnasirpal II and Shalmaneser III. The arts in particular experienced a tremendous resurgence. Literary works continued to be written in Assyrian and were seldom of great importance. The literature that had been taken over from Babylonia was further developed with new writings, although one can rarely distinguish between works written in Assyria and works written in Babylonia. In religion, the official cults of Ashur and Ninurta continued, while the religion of the common people went its separate way.

In Babylonia not much was left of the feudal structure; the large landed estates almost everywhere fell prey to the inroads of the Aramaeans, who were at first half nomadic. The leaders of their tribes and clans slowly replaced the former landlords. Agriculture on a large scale was no longer possible except on the outskirts of metropolitan areas. The predominance of the Babylonian schools for scribes may have prevented the emergence of an Aramaean literature. In any case, the Aramaeans seem to have been absorbed into the Babylonian culture. The religious cults in the cities remained essentially the same. The Babylonian empire was slowly reduced to poverty, except perhaps in some of the cities.

In 764, after an epidemic, the *Erra* epic, the myth of *Erra* (the god of war and pestilence), was written by Kabti-ilani-Marduk. He invented an original plot, which diverged considerably from the old myths; long discourses of the gods involved in the action form the most important part of the epic. There is a passage in the epic claiming that the text was divinely revealed to the poet during a dream.

THE NEO-ASSYRIAN EMPIRE (746–609)

For no other period of Assyrian history is there an abundance of sources comparable to those available for the interval from roughly 745 to 640. Aside from the large number of royal inscriptions, about 2,400 letters, most of them more or less fragmentary, have been published. Usually the senders and recipients of these letters are the king and high government officials. Among them are reports from royal agents about foreign affairs and letters about cultic matters. Treaties, oracles, queries to the sun god about political matters, and prayers of or for kings contain a great deal of additional information. Last but certainly not least are paintings and wall reliefs, which are often very informative.

Tiglath-pileser III and Shalmaneser V. The decline of Assyrian power after 780 was notable; Syria and considerable lands in the north were lost. A military coup deposed King Ashur-nirari V and raised a general to the throne. Under the name of Tiglath-pileser III (746–727), he brought the empire to its greatest expanse. He reduced the size of the provinces in order to break the partial independence of the governors. He also invalidated the tax privileges of cities such as Ashur and Harran in order to distribute the tax load more evenly over the entire realm. Military equipment was improved substantially. In 746 he went to Babylonia to aid Nabu-naṣir (747–734) in his fight against Aramaean tribes. Tiglath-pileser defeated the Aramaeans and then made visits to the large cities of Babylonia. There he tried to secure the support of the priesthood by patronizing their building projects. Babylonia retained its independence.

Changes
in social
structure

His next undertaking was to check Urartu. His campaigns in Azerbaijan were designed to drive a wedge between Urartu and the Medes. In 743 he went to Syria, defeating there an army of Urartu. The Syrian city of Arpad, which had formed an alliance with Urartu, did not surrender so easily. It took Tiglath-pileser three years of siege to conquer Arpad, whereupon he massacred the inhabitants and destroyed the city. In 738 a new coalition formed against Assyria under the leadership of Sam'al (modern Zincirli) in northern Syria. It was defeated, and all the princes from Damascus to eastern Anatolia were forced to pay tribute. Another campaign in 735, this time directed against Urartu itself, was only partly successful. In 734 Tiglath-pileser invaded southern Syria and the Philistine territories in Palestine, going as far as the Egyptian border. Damascus and Israel tried to organize resistance against him, seeking to bring Judah into their alliance. Ahaz of Judah, however, asked Tiglath-pileser for help. In 733 Tiglath-pileser devastated Israel and forced it to surrender large territories. In 732 he advanced upon Damascus, first devastating the gardens outside the city and then conquering the capital and killing the king, whom he replaced with a governor. The queen of southern Arabia, Samsil, was now obliged to pay tribute, being permitted in return to use the harbour of the city of Gaza, which was in Assyrian hands.

The death of King Nabonassar of Babylonia caused a chaotic situation to develop there, and the Aramaean Ukin-zer crowned himself king. In 731 Tiglath-pileser fought and beat him and his allies, but he did not capture Ukin-zer until 729. This time he did not appoint a new king for Babylonia but assumed the crown himself under the name Pulu (Pul in the Old Testament). In his old age he abstained from further campaigning, devoting himself

to the improvement of his capital, Kalakh. He rebuilt the palace of Shalmaneser III, filled it with treasures from his wars, and decorated the walls with bas-reliefs. The latter were almost all of warlike character, as if designed to intimidate the onlooker with their presentation of gruesome executions. These pictorial narratives on slabs, sometimes painted, have also been found in Syria, at the sites of several provincial capitals of ancient Assyria.

Tiglath-pileser was succeeded by his son Shalmaneser V (726-722), who continued the policy of his father. As king of Babylonia, he called himself Ululai. Almost nothing is known about his enterprises, since his successor destroyed all his inscriptions. The Old Testament relates that he marched against Hoshea of Israel in 724 after Hoshea had rebelled. He was probably assassinated during the long siege of Samaria. His successor maintained that the god Ashur had withdrawn his support of Shalmaneser V for acts of disrespect.

Sargon II (721-705) and Marduk-apal-iddina of Babylonia. It was probably a younger brother of Shalmaneser who ascended the throne of Assyria in 721. Assuming the old name of Sharru-kin (Sargon in the Bible), meaning "Legitimate King," he assured himself of the support of the priesthood and the merchant class by restoring privileges they had lost, particularly the tax exemptions of the great temples. The change of sovereign in Assyria triggered another crisis in Babylonia. An Aramaean prince from the south, Marduk-apal-iddina II (the biblical Merodach-Baladan), seized power in Babylon in 721 and was able to retain it until 710 with the help of Humbanigash I of Elam. A first attempt by Sargon to recover Babylonia miscarried when Elam defeated him in 721. During the same year the protracted siege of Samaria was brought to a close. The Samarian upper class was deported, and Israel became an

Wars against Urartu

Campaigns of Sargon II

From W. Shepherd, *Historical Atlas*, Harper & Row, Publishers (Barnes & Noble Books), New York, revision copyright © 1964 by Barnes & Noble, Inc.



The Assyrian empire, 858-627 BC.

Assyrian province. Samaria was repopulated with Syrians and Babylonians. Judah remained independent by paying tribute. In 720 Sargon squelched a rebellion in Syria that had been supported by Egypt. Then he defeated both Hanunu of Gaza and an Egyptian army near the Egyptian border. In 717 and 716 he campaigned in northern Syria, making the hitherto independent state of Carchemish one of his provinces. He also went to Cilicia in an effort to prevent further encroachments of the Phrygians under King Midas (Assyrian: Mitā).

In order to protect his ally, the state of Mannai, in Azerbaijan, Sargon embarked on a campaign in Iran in 719 and incorporated parts of Media as provinces of his empire; however, in 716 another war became necessary. At the same time, he was busy preparing a major attack against Urartu. Under the leadership of the crown prince Sennacherib, armies of agents infiltrated Urartu, which was also threatened from the north by the Cimmerians. Many of their messages and reports have been preserved. The longest inscription ever composed by the Assyrians about a year's enterprise (430 very long lines) is dedicated to this Urartu campaign of 714. Phrased in the style of a first report to the god Ashur, it is interspersed with stirring descriptions of natural scenery. The strong points of Urartu must have been well fortified. Sargon tried to avoid them by going through the province of Mannai and attacking the Median principalities on the eastern side of Lake Urmia. In the meantime, hoping to surprise the Assyrian troops, Rusa of Urartu had closed the narrow pass lying between Lake Urmia and Sahand Mount. Sargon, anticipating this, led a small band of cavalry in a surprise charge that developed into a great victory for the Assyrians. Rusa fled and died. The Assyrians pushed forward, destroying all the cities, fortifications, and even irrigation works of Urartu. They did not conquer Tushpa (the capital) but took possession of the mountain city of Muṣaṣir. The spoils were immense. The following years saw only small campaigns in Media and eastern Anatolia and against Ashdod, in Palestine. King Midas of Phrygia and some cities on Cyprus were quite ready to pay tribute.

Sargon was now free to settle accounts with Marduk-apal-iddina of Babylonia. Abandoned by his ally Shutrūk-Nahhunte II of Elam, Marduk-apal-iddina found it best to flee, first to his native land on the Persian Gulf and later to Elam. Because the Aramaean prince had made himself very unpopular with his subjects, Sargon was hailed as the liberator of Babylonia. He complied with the wishes of the priesthood and at the same time put down the Aramaean nobility. He was satisfied with the modest title of governor of Babylonia.

At first Sargon resided in Kalakh, but he then decided to found an entirely new capital north of Nineveh. He called the city Dur-Sharrukin—"Sargonsburg" (modern Khorsabad, Iraq). He erected his palace on a high terrace in the northeastern part of the city. The temples of the main gods, smaller in size, were built within the palatial rectangle, which was surrounded by a special wall. This arrangement enabled Sargon to supervise the priests better than had been possible in the old, large temple complexes. One consequence of this design was that the figure of the king pushed the gods somewhat into the background, thereby gaining in importance. Desiring that his palace match the vastness of his empire, Sargon planned it in monumental dimensions. Stone reliefs of two winged bulls with human heads flanked the entrance; they were much larger than anything comparable built before. The walls were decorated with long rows of bas-reliefs showing scenes of war and festive processions. A comparison with a well-executed stela of the Babylonian king Marduk-apal-iddina shows that the fine arts of Assyria had far surpassed those of Babylonia. Sargon never completed his capital, though from 713 to 705 bc tens of thousands of labourers and hundreds of artisans worked on the great city. Yet, with the exception of some magnificent buildings for public officials, only a few durable edifices were completed in the residential section. In 705, in a campaign in northwestern Iran, Sargon was ambushed and killed. His corpse remained unburied, to be devoured by birds of prey. Sargon's son Sennacherib, who had quarreled with

his father, was inclined to believe with the priests that his death was a punishment from the neglected gods of the ancient capitals.

Sennacherib. Sennacherib (Assyrian: Sin-ahhe-eriba; 704–681) was well prepared for his position as sovereign. With him Assyria acquired an exceptionally clever and gifted, though often extravagant, ruler. His father, interestingly enough, is not mentioned in any of his many inscriptions. He left the new city of Dur-Sharrukin at once and resided in Ashur for a few years, until in 701 he made Nineveh his capital.

Sennacherib had considerable difficulties with Babylonia. In 703 Marduk-apal-iddina again crowned himself king with the aid of Elam, proceeding at once to ally himself with other enemies of Assyria. After nine months he was forced to withdraw when Sennacherib defeated a coalition army consisting of Babylonians, Aramaeans, and Elamites. The new puppet king of Babylonia was Bel-ibni (702–700), who had been raised in Assyria.

In 702 Sennacherib launched a raid into western Iran. In 701 there followed his most famous campaign, against Syria and Palestine, with the purpose of gaining control over the main road from Syria to Egypt in preparation for later campaigns against Egypt itself. When Sennacherib's army approached, Sidon immediately expelled its ruler, Luli, who was hostile to Assyria. The other allies either surrendered or were defeated. An Egyptian army was defeated at Eltekeh in Judah. Sennacherib laid siege to Jerusalem, and the king of Judah, Hezekiah, was called upon to surrender, but he did not comply. An Assyrian officer tried to incite the people of Jerusalem against Hezekiah, but his efforts failed. In view of the difficulty of surrounding a mountain stronghold such as Jerusalem, and of the minor importance of this town for the main purpose of the campaign, Sennacherib cut short the attack and left Palestine with his army, which according to the Old Testament (2 Kings 19:35) had been decimated by an epidemic. The number of Assyrian dead is reported to have risen to 185,000. Nevertheless, Hezekiah is reported to have paid tribute to Sennacherib on at least one occasion.

Bel-ibni of Babylonia seceded from the union with Assyria in 700. Sennacherib moved quickly, defeating Bel-ibni and replacing him with Sennacherib's oldest son, Ashur-nadin-shumi. The next few years were relatively peaceful. Sennacherib used this time to prepare a decisive attack against Elam, which time and again had supported Babylonian rebellions. The overland route to Elam had been cut off and fortified by the Elamites. Sennacherib had ships built in Syria and at Nineveh. The ships from Syria were moved on rollers from the Euphrates to the Tigris. The fleet sailed downstream and was quite successful in the lagoons of the Persian Gulf and along the southern coastline of Elam. The Elamites launched a counteroffensive by land, occupying Babylonia and putting a man of their choice on the throne. Not until 693 were the Assyrians again able to fight their way through to the north. Finally, in 689, Sennacherib had his revenge. Babylon was conquered and completely destroyed, the temples plundered and leveled. The waters of the Arakhtu Canal were diverted over the ruins, and the inner city remained almost totally uninhabited for eight years. Even many Assyrians were indignant at this, believing that the Babylonian god Marduk must be grievously offended at the destruction of his temple and the carrying off of his image. Marduk was also an Assyrian deity, to whom many Assyrians turned in time of need. A political-theological propaganda campaign was launched to explain to the people that what had taken place was in accord with the wish of most of the gods. A story was written in which Marduk, because of a transgression, was captured and brought before a tribunal. Only a part of the commentary to this botched piece of literature is extant. Even the great poem of the creation of the world, the *Enuma elish*, was altered: the god Marduk was replaced by the god Ashur. Sennacherib's boundless energies brought no gain to his empire, however, and probably weakened it. The tenacity of this king can be seen in his building projects; for example, when Nineveh needed water for irrigation, Sennacherib had his engineers

The
siege of
Jerusalem

Sargon's
palace

divert the waters of a tributary of the Great Zab River. The canal had to cross a valley at Jerwan. An aqueduct was constructed, consisting of about two million blocks of limestone, with five huge, pointed archways over the brook in the valley. The bed of the canal on the aqueduct was sealed with cement containing magnesium. Parts of this aqueduct are still standing today. Sennacherib wrote of these and other technological accomplishments in minute detail, with illustrations.

Senna-
cherib's
palace in
Nineveh

Sennacherib built a huge palace in Nineveh, adorned with reliefs, some of them depicting the transport of colossal bull statues by water and by land. Many of the rooms were decorated with pictorial narratives in bas-relief telling of war and of building activities. Considerable advances can be noted in artistic execution, particularly in the portrayal of landscapes and animals. Outstanding are the depictions of the battles in the lagoons, the life in the military camps, and the deportations.

In 681 BC there was a rebellion. Sennacherib was assassinated by one or two of his sons in the temple of the god Ninurta at Kalakh. This god, along with the god Marduk, had been badly treated by Sennacherib, and the event was widely regarded as punishment of divine origin.

Esarhaddon. Ignoring the claims of his older brothers, an imperial council appointed Esarhaddon (Ashur-ah-iddina; 680–669) as Sennacherib's successor. The choice is all the more difficult to explain in that Esarhaddon, unlike his father, was friendly toward the Babylonians. It can be assumed that his energetic and designing mother, Zakutu (Naqia), who came from Syria or Judah, used all her influence on his behalf to override the national party of Assyria. The theory that he was a partner in plotting the murder of his father is rather improbable; at any rate, he was able to procure the loyalty of his father's army. His brothers had to flee to Urartu. In his inscriptions, Esarhaddon always mentions both his father and grandfather.

Defining the destruction of Babylon explicitly as punishment by the god Marduk, the new king soon ordered the reconstruction of the city. He referred to himself only as governor of Babylonia and through his policies obtained the support of the cities of Babylonia. At the beginning of his reign the Aramaean tribes were still allied with Elam against him, but Urtaku of Elam (675–664) signed a peace treaty and freed him for campaigning elsewhere. In 679 he stationed a garrison at the Egyptian border, because Egypt, under the Ethiopian king Taharqa, was planning to intervene in Syria. He put down with great severity a rebellion of the combined forces of Sidon, Tyre, and other Syrian cities. The time was ripe to attack Egypt, which was suffering under the rule of the Ethiopians and was by no means a united country. Esarhaddon's first attempt in 674–673 miscarried. In 671 BC, however, his forces took Memphis, the Egyptian capital. Assyrian consultants were assigned to assist the princes of the 22 provinces, their main duty being the collection of tribute.

Occasional threats came from the mountainous border regions of eastern Anatolia and Iran. Pushed forward by the Scythians, the Cimmerians in northern Iran and Transcaucasia tried to gain a foothold in Syria and western Iran. Esarhaddon allied himself with the Scythian king Partatua by giving him one of his daughters in marriage. In so doing he checked the movement of the Cimmerians. Nevertheless, the apprehensions of Esarhaddon can be seen in his many offerings, supplications, and requests to the sun god. These were concerned less with his own enterprises than with the plans of enemies and vassals and the reliability of civil servants. The priestesses of Ishtar had to reassure Esarhaddon constantly by calling out to him, "Do not be afraid." Previous kings, as far as is known, had never needed this kind of encouragement.

At home Esarhaddon was faced with serious difficulties from factions in the court. His oldest son had died early. The national party suspected his second son, Shamash-shum-ukin, of being too friendly with the Babylonians; he may also have been considered unequal to the task of kingship. His third son, Ashurbanipal, was given the succession in 672, Shamash-shum-ukin remaining crown prince of Babylonia. This arrangement caused much dissension, and some farsighted civil servants warned of di-

Domestic
problems

sastrous effects. Nevertheless, the Assyrian nobles, priests, and city leaders were sworn to just such an adjustment of the royal line; even the vassal princes had to take very detailed oaths of allegiance to Ashurbanipal, with many curses against perjurers.

Another matter of deep concern for Esarhaddon was his failing health. He regarded eclipses of the moon as particularly alarming omens, and, in order to prevent a fatal illness from striking him at these times, he had substitute kings chosen who ruled during the three eclipses that occurred during his 12-year reign. The replacement kings died or were put to death after their brief term of office. During his off-terms Esarhaddon called himself "Mister Peasant." This practice implied that the gods could not distinguish between the real king and a false one—quite contrary to the usual assumptions of the religion.

Esarhaddon enlarged and improved the temples in both Assyria and Babylonia. He also constructed a palace in Kalakh, using many of the picture slabs of Tiglath-pileser III. The works that remain are not on the level of those of either his predecessors or of Ashurbanipal. He died while on an expedition to put down a revolt in Egypt.

Ashurbanipal (668–627) and Shamash-shum-ukin (668–648). Although the death of his father occurred far from home, Ashurbanipal assumed the kingship as planned. He may have owed his fortunes to the intercession of his grandmother Zakutu, who had recognized his superior capacities. He tells of his diversified education by the priests and his training in armour-making as well as in other military arts. He may have been the only king in Assyria with a scholarly background. As crown prince he also had studied the administration of the vast empire. The record notes that the gods granted him a record harvest during the first year of his reign. There were also good crops in subsequent years. During these first years he also was successful in foreign policy, and his relationship with his brother in Babylonia was good.

In 668 he put down a rebellion in Egypt and drove out King Taharqa, but in 664 the nephew of Taharqa, Tanutamon, gathered forces for a new rebellion. Ashurbanipal went to Egypt, pursuing the Ethiopian prince far into the south. His decisive victory moved Tyre and other parts of the empire to resume regular payments of tribute. Ashurbanipal installed Psamtik (Greek: Psammetichos) as prince over the Egyptian region of Sais. In 656 Psamtik dislodged the Assyrian garrisons with the aid of Carian and Ionian mercenaries, making Egypt again independent. Ashurbanipal did not attempt to reconquer it. A former ally of Assyria, Gyges of Lydia, had aided Psamtik in his rebellion. In return, Assyria did not help Gyges when he was attacked by the Cimmerians. Gyges lost his throne and his life. His son Ardys decided that the payment of tribute to Assyria was a lesser evil than conquest by the Cimmerians.

Graver difficulties loomed in southern Babylonia, which was attacked by Elam in 664. Another attack came in 653, whereupon Ashurbanipal sent a large army that decisively defeated the Elamites. Their king was killed, and some of the Elamite states were encouraged to secede. Elam was no longer strong enough to assume an active part on the international scene. This victory had serious consequences for Babylonia. Shamash-shum-ukin had grown weary of being patronized by his domineering brother. He formed a secret alliance in 656 with the Iranians, Elamites, Aramaeans, Arabs, and Egyptians, directed against Ashurbanipal. The withdrawal of defeated Elam from this alliance was probably the reason for a premature attack by Shamash-shum-ukin at the end of the year 652, without waiting for the promised assistance from Egypt. Ashurbanipal, taken by surprise, soon pulled his troops together. The Babylonian army was defeated, and Shamash-shum-ukin was surrounded in his fortified city of Babylon. His allies were not able to hold their own against the Assyrians. Reinforcements of Arabian camel troops also were defeated. The city of Babylon was under siege for three years. It fell in 648 amid scenes of horrible carnage, Shamash-shum-ukin dying in his burning palace.

After 648 the Assyrians made a few punitive attacks on the Arabs, breaking the forward thrust of the Arab tribes

War
with
Elam

for a long time to come. The main objective of the Assyrians, however, was a final settlement of their relations with Elam. The refusal of Elam in 647 to extradite an Aramaean prince was used as pretext for a new attack that drove deep into its territory. The assault on the solidly fortified capital of Susa followed, probably in 646. The Assyrians destroyed the city, including its temples and palaces. Vast spoils were taken. As usual, the upper classes of the land were exiled to Assyria and other parts of the empire, and Elam became an Assyrian province. Assyria had now extended its domain to southwestern Iran. Cyrus I of Persia sent tribute and hostages to Nineveh, hoping perhaps to secure protection for his borders with Media. Little is known about the last years of Ashurbanipal's reign.

Ashurbanipal left more inscriptions than any of his predecessors. His campaigns were not always recorded in chronological order but clustered in groups according to their purpose. The accounts were highly subjective. One of his most remarkable accomplishments was the founding of the great palace library in Nineveh (modern Kuyunjik), which is today one of the most important sources for the study of ancient Mesopotamia. The king himself supervised its construction. Important works were kept in more than one copy, some intended for the king's personal use. The work of arranging and cataloging drew upon the experience of centuries in the management of collections in huge temple archives such as the one in Ashur. In his inscriptions Ashurbanipal tells of becoming an enthusiastic hunter of big game, acquiring a taste for it during a fight with marauding lions. In his palace at Nineveh the long rows of hunting scenes show what a masterful artist can accomplish in bas-relief; with these reliefs Assyrian art reached its peak. In the series depicting his wars, particularly the wars fought in Elam, the scenes are overloaded with human figures. Those portraying the battles with the Arabian camel troops are magnificent in execution.

The policy of resettlement

One reason for the durability of the Assyrian empire was the practice of deporting large numbers of people from conquered areas and resettling others in their place. This kept many of the conquered nationalities from regaining their power. Equally important was the installation in conquered areas of a highly developed civil service under the leadership of trained officers. The highest ranking civil servant carried the title of *tartān*, a Hurrian word. The *tartāns* also represented the king during his absence. In descending rank were the palace overseer, the main cupbearer, the palace administrator, and the governor of Assyria. The generals often held high official positions, particularly in the provinces. The civil service numbered about 100,000, many of them former inhabitants of subjugated provinces. Prisoners became slaves but were later often freed.

No laws are known for the empire, although documents point to the existence of rules and standards for justice. Those who broke contracts were subject to severe penalties, even in cases of minor importance: the sacrifice of a son or the eating of a pound of wool and drinking of a great deal of water afterward, which led to a painful death. The position of women was inferior, except for the queen and some priestesses.

As yet there are no detailed studies of the economic situation during this period. The landed nobility still played an important role, in conjunction with the merchants in the cities. The large increase in the supply of precious metals—received as tribute or taken as spoils—did not disrupt economic stability in many regions. Stimulated by the patronage of the kings and the great temples, the arts and crafts flourished during this period. The policy of resettling Aramaeans and other conquered peoples in Assyria brought many talented artists and artisans into Assyrian cities, where they introduced new styles and techniques. High-ranking provincial civil servants, who were often very powerful, saw to it that the provincial capitals also benefited from this economic and cultural growth.

Harran became the most important city in the western part of the empire; in the neighbouring settlement of Huzirina (modern Sultantepe, in northern Syria), the remains of an important library have been discovered. Very

few Aramaic texts from this period have been found; the climate of Mesopotamia is not conducive to the preservation of the papyrus and parchment on which these texts were written. There is no evidence that a literary tradition existed in any of the other languages spoken within the borders of the Assyrian empire at this time, except in peripheral areas of Syria and Palestine.

Culturally and economically, Babylonia lagged behind Assyria in this period. The wars with Assyria—particularly the catastrophic defeats of 689 and 648—together with many smaller tribal wars disrupted trade and agricultural production. The great Babylonian temples fared best during this period, since they continued to enjoy the patronage of the Assyrian monarchs. Only a few documents from the temples have been preserved, however. There is evidence that the scribal schools continued to operate, and "Sumerian" inscriptions were even composed for Shamash-shum-ukin. In comparison with the Assyrian developments, the pictorial arts were neglected, and Babylonian artists may have found work in Assyria.

During this period people began to use the names of ancestors as a kind of family name; this increase in family consciousness is probably an indication that the number of old families was growing smaller. By this time the process of "Aramaicization" had reached even the oldest cities of Babylonia and Assyria.

Apparently this era was not very fruitful for literature either in Babylonia or in Assyria. In Assyria numerous royal inscriptions, some as long as 1,300 lines, were among the most important texts; some of them were diverse in content and well composed. Most of the hymns and prayers were written in the traditional style. Many oracles, often of unusual content, were proclaimed in the Assyrian dialect, most often by the priestesses of the goddess Ishtar of Arbela. In Assyria as in Babylonia, the beginnings of a real historical literature are observed; most of the authors have remained anonymous up to the present.

Assyrian royal inscriptions

The many gods of the tradition were worshiped in Babylonia and Assyria in large and small temples, as in earlier times. Very detailed rituals regulated the sacrifices, and the interpretations of the ritual performances in the cultic commentaries were rather different and sometimes very strange.

On some of the temple towers (ziggurats), astronomical observatories were installed. The earliest of these may have been the observatory of the Ninurta temple at Kalakh in Assyria, which dates back to the 9th century BC; it was destroyed with the city in 612. The most important observatory in Babylonia from about 580 was situated on the ziggurat Etemenanki, a temple of Marduk in Babylon. In Assyria the observation of the Sun, Moon, and stars had already reached a rather high level; the periodic recurrence of eclipses was established. After 600, astronomical observation and calculations developed steadily, and they reached their high point after 500, when Babylonian and Greek astronomers began their fruitful collaboration. Incomplete astronomical diaries, beginning in 652 and covering some 600 years, have been preserved.

Decline of the Assyrian empire. Few historical sources remain for the last 30 years of the Assyrian empire. There are no extant inscriptions of Ashurbanipal after 640 BC, and the few surviving inscriptions of his successors contain only vague allusions to political matters. In Babylonia the silence is almost total until 625 BC, when the chronicles resume. The rapid downfall of the Assyrian empire was formerly attributed to military defeat, although it was never clear how the Medes and the Babylonians alone could have accomplished this. More recent work has established that after 635 a civil war occurred, weakening the empire so that it could no longer stand up against a foreign enemy. Ashurbanipal had twin sons. Ashur-etel-ilani was appointed successor to the throne, but his twin brother Sin-shar-ishkun did not recognize him. The fight between them and their supporters forced the old king to withdraw to Harran, in 632 at the latest, perhaps ruling from there over the western part of the empire until his death in 627. Ashur-etel-ilani governed in Assyria from about 633, but a general, Sin-shum-lisher, soon rebelled against him and proclaimed himself counter-king. Some

Internal struggles in the last years of the empire

years later (629?) Sin-shar-ishkun finally succeeded in obtaining the kingship. In Babylonian documents dates can be found for all three kings. To add to the confusion, until 626 there are also dates of Ashurbanipal and a king named Kandalanu. In 626 the Chaldean Nabopolassar (Nabu-apal-usur) revolted from Uruk and occupied Babylon. There were several changes in government. King Ashur-etel-ilani was forced to withdraw to the west, where he died sometime after 625.

About the year 626 the Scythians laid waste to Syria and Palestine. In 625 the Medes became united under Cyaxares and began to conquer the Iranian provinces of Assyria. One chronicle relates of wars between Sin-shar-ishkun and Nabopolassar in Babylonia in 625–623. It was not long until the Assyrians were driven out of Babylonia. In 616 the Medes struck against Nineveh, but, according to the Greek historian Herodotus, were driven back by the Scythians. In 615, however, the Medes conquered Arapkha (Kirkük), and in 614 they took the old capital of Ashur, looting and destroying the city. Now Cyaxares and Nabopolassar made an alliance for the purpose of dividing Assyria. In 612 Kalakh and Nineveh succumbed to the superior strength of the allies. The revenge taken on the Assyrians was terrible: 200 years later Xenophon found the country still sparsely populated.

Sin-shar-ishkun, king of Assyria, found death in his burning palace. The commander of the Assyrian army in the west crowned himself king in the city of Harran, assuming the name of the founder of the empire, Ashur-uballit II (611–609 BC). Ashur-uballit had to face both the Babylonians and the Medes. They conquered Harran in 610, without, however, destroying the city completely. In 609 the remaining Assyrian troops had to capitulate. With this event Assyria disappeared from history. The great empires that succeeded it learned a great deal from the hated Assyrians, both in the arts and in the organization of their states.

THE NEO-BABYLONIAN EMPIRE

The Chaldeans, who inhabited the coastal area near the Persian Gulf, had never been entirely pacified by the Assyrians. About 630 Nabopolassar became king of the Chaldeans. In 626 he forced the Assyrians out of Uruk and crowned himself king of Babylonia. He took part in the wars aimed at the destruction of Assyria. At the same time, he began to restore the dilapidated network of canals in the cities of Babylonia, particularly those in Babylon itself. He fought against the Assyrian Ashur-uballit II and then against Egypt, his successes alternating with misfortunes. In 605 Nabopolassar died in Babylon.

Nebuchadnezzar II. Nabopolassar had named his oldest son, Nabu-kudurri-ušur, after the famous king of the second dynasty of Isin, trained him carefully for his prospective kingship, and shared responsibility with him. When the father died in 605, Nebuchadnezzar was with his army in Syria; he had just crushed the Egyptians near Carchemish in a cruel, bloody battle and pursued them into the south. On receiving the news of his father's death, Nebuchadnezzar returned immediately to Babylon. In his numerous building inscriptions he tells but rarely of his many wars; most of them end with prayers. The Babylonian chronicle is extant only for the years 605–594, and not much is known from other sources about the later years of this famous king. He went very often to Syria and Palestine, at first to drive out the Egyptians. In 604 he took the Philistine city of Ashkelon. In 601 he tried to push forward into Egypt but was forced to pull back after a bloody, undecided battle and to regroup his army in Babylonia. After smaller incursions against the Arabs of Syria, he attacked Palestine at the end of 598. King Jehoiakim of Judah had rebelled, counting on help from Egypt. According to the chronicle, Jerusalem was taken on March 16, 597. Jehoiakim had died during the siege, and his son, King Jehoiachin, together with at least 3,000 Jews, was led into exile in Babylonia. They were treated well there, according to the documents. Zedekiah was appointed the new king. In 596, when danger threatened from the east, Nebuchadnezzar marched to the Tigris River and induced the enemy to withdraw. After a revolt

in Babylonia had been crushed with much bloodshed, there were other campaigns in the west.

According to the Old Testament, Judah rebelled again in 589, and Jerusalem was placed under siege. The city fell in 587/586 and was completely destroyed. Many thousands of Jews were forced into "Babylonian exile," and their country was reduced to a province of the Babylonian empire. The revolt had been caused by an Egyptian invasion that pushed as far as Sidon. Nebuchadnezzar laid siege to Tyre for 13 years without taking the city, because there was no fleet at his disposal. In 568/567 he attacked Egypt, again without much success, but from that time on the Egyptians refrained from further attacks on Palestine. Nebuchadnezzar lived at peace with Media throughout his reign and acted as a mediator after the Median-Lyidian war of 590–585.

The Babylonian empire under Nebuchadnezzar extended to the Egyptian border. It had a well-functioning administrative system. Though he had to collect extremely high taxes and tributes in order to maintain his armies and carry out his building projects, Nebuchadnezzar made Babylonia one of the richest lands in western Asia—the more astonishing because it had been rather poor when it was ruled by the Assyrians. Babylon was the largest city of the "civilized world." Nebuchadnezzar maintained the existing canal systems and built many supplementary canals, making the land even more fertile. Trade and commerce flourished during his reign.

Nebuchadnezzar's building activities surpassed those of most of the Assyrian kings. He fortified the old double walls of Babylon, adding another triple wall outside the old wall. In addition, he erected another wall, the Median Wall, north of the city between the Euphrates and the Tigris rivers. According to Greek estimates, the Median Wall may have been about 100 feet high. He enlarged the old palace and added many wings, so that hundreds of rooms with large inner courts were now at the disposal of the central offices of the empire. Colourful glazed-tile bas-reliefs decorated the walls. Terrace gardens, called the Hanging Gardens in later accounts, were added. Hundreds of thousands of workers must have been required for these projects. The temples were objects of special concern. He devoted himself first and foremost to the completion of Etemenanki, the "Tower of Babel." Construction of this building began in the time of Nebuchadnezzar I, about 1110. It stood as a "building ruin" until the reign of Esarhaddon of Assyria, who resumed building about 680 but did not finish. Nebuchadnezzar II was able to complete the whole building. The mean dimensions of Etemenanki are to be found in the Esagila Tablet, which has been known since the late 19th century. Its base measured about 300 feet on each side, and it was 300 feet in height. There were five terraced gradations surmounted by a temple, the whole tower being about twice the height of those of other temples. The wide street used for processions led along the eastern side by the inner city walls and crossed at the enormous Ishtar Gate with its world-renowned bas-relief tiles. Nebuchadnezzar also built many smaller temples throughout the country.

The last kings of Babylonia. Awil-Marduk (called Evil-Merodach in the Old Testament; 561–560), the son of Nebuchadnezzar, was unable to win the support of the priests of Marduk. His reign did not last long, and he was soon eliminated. His brother-in-law and successor, Nergal-shar-ušur (called Neriglissar in classical sources; 559–556), was a general who undertook a campaign in 557 into the "rough" Cilician land, which may have been under the control of the Medes. His land forces were assisted by a fleet. His still-minor son Labashi-Marduk was murdered not long after that, allegedly because he was not suitable for his job.

The next king was the Aramean Nabonidus (Nabunaid; 556–539) from Harran, one of the most interesting and enigmatic figures of ancient times. His mother, Addagoppe, was a priestess of the god Sin in Harran; she came to Babylon and managed to secure responsible offices for her son at court. The god of the moon rewarded her piety with a long life—she lived to be 103—and she was buried in Harran with all the honours of

Jerusalem
destroyed

The
Tower
of Babel

Assyrians
driven
from
Babylonia

Campaigns
in Syria
and
Palestine

a queen in 547. It is not clear which powerful faction in Babylon supported the kingship of Nabonidus; it may have been one opposing the priests of Marduk, who had become extremely powerful. Nabonidus raided Cilicia in 555 and secured the surrender of Harran, which had been ruled by the Medes. He concluded a treaty of defense with Astyages of Media against the Persians, who had become a growing threat since 559 under their king Cyrus II. He also devoted himself to the renovation of many temples, taking an especially keen interest in old inscriptions. He gave preference to his god Sin and had powerful enemies in the priesthood of the Marduk temple. Modern excavators have found fragments of propaganda poems written against Nabonidus and also in support of him. Both traditions continued in Judaism.

Internal difficulties and the recognition that the narrow strip of land from the Persian Gulf to Syria could not be defended against a major attack from the east induced Nabonidus to leave Babylonia around 552 and to reside in Taima (Taymā') in northern Arabia. There he organized an Arabian province with the assistance of Jewish mercenaries. His viceroy in Babylonia was his son Bel-shar-ušur, the Belshazzar of the Book of Daniel in the Bible. Cyrus turned this to his own advantage by annexing Media in 550. Nabonidus, in turn, allied himself with Croesus of Lydia in order to fight Cyrus. Yet, when Cyrus attacked Lydia and annexed it in 546, Nabonidus was not able to help Croesus. Cyrus bode his time. In 542 Nabonidus returned to Babylonia, where his son had been able to maintain good order in external matters but had not overcome a growing internal opposition to his father. Consequently, Nabonidus' career after his return was short-lived, though he tried hard to regain the support of the Babylonians. He appointed his daughter to be high priestess of the god Sin in Ur, thus returning to the Sumerian-Old Babylonian religious tradition. The priests of Marduk looked to Cyrus, hoping to have better relations with him than with Nabonidus; they promised Cyrus the surrender of Babylon without a fight if he would grant them their privileges in return. In 539 Cyrus attacked northern Babylonia with a large army, defeating Nabonidus, and entered the city of Babylon without a battle. The other cities did not offer any resistance either. Nabonidus surrendered, receiving a small territory in eastern Iran. Tradition has confused him with his great predecessor Nebuchadrezzar II. The Bible refers to him as Nebuchadrezzar in the Book of Daniel.

Babylonia's peaceful submission to Cyrus saved it from the fate of Assyria. It became a territory under the Persian crown but kept its cultural autonomy. Even the racially mixed western part of the Babylonian empire submitted without resistance.

By 620 the Babylonians had grown tired of Assyrian rule. They were also weary of internal struggle. They were easily persuaded to submit to the order of the Chaldean kings. The result was a surprisingly rapid social and economic consolidation, helped along by the fact that after the fall of Assyria no external enemy threatened Babylonia for more than 60 years. In the cities the temples were an important part of the economy, having vast benefices at their disposal. The business class regained its strength, not only in the trades and commerce but also in the management of agriculture in the metropolitan areas. Livestock breeding—sheep, goats, beef cattle, and horses—flourished, as did poultry farming. The cultivation of corn, dates, and vegetables grew in importance. Much was done to improve communications, both by water and land, with the western provinces of the empire. The collapse of the Assyrian empire had the consequence that many trade arteries were rerouted through Babylonia. Another result of the collapse was that the city of Babylon became a world centre.

The immense amount of documentary material and correspondence that has survived has not yet been fully analyzed. No new system of law or administration seems to have developed during that time. The Babylonian dialect gradually became Aramaicized; it was still written primarily on clay tablets that often bore added material in Aramaic lettering. Parchment and papyrus documents have not survived. In contrast to advances in other fields, there is no evidence of much artistic creativity. Aside from

some of the inscriptions of the kings, especially Nabonidus, which were not comparable from a literary standpoint with those of the Assyrians, the main efforts were devoted to the rewriting of old texts. In the fine arts, only a few monuments have any suggestion of new tendencies.

MESOPOTAMIA UNDER THE PERSIANS

Cyrus II, the founder of the Achaemenian Empire, united Babylonia with his country in a personal union, assuming the title of "King of Babylonia, King of the Lands." His son Cambyses was appointed vice-king and resided in Sippar. The Persians relied on the support of the priests and the business class in the cities. In a Babylonian inscription, Cyrus relates with pride his peaceful, bloodless conquest of the city of Babylon. At the same time, he speaks of Marduk as the king of gods. His moderation and restraint were rewarded: Babylonia became the richest province of his empire. There is no indication of any national rebellion in Babylonia under Cyrus and Cambyses (529–522). That there must have been an accumulation of discontent became clear at the ascension to the throne of Darius I (522–486), when a usurper seized the throne of Babylonia under the name of Nebuchadrezzar (III) only to lose both the throne and his life after 10 weeks. Darius waived any punitive action. He had to take more drastic measures in 521, when a new Nebuchadrezzar incited another rebellion. This usurper's reign lasted two months. Executions and plundering followed; Darius ordered that the inner walls of Babylon be demolished, and he reformed the organization of the state. Babylon, however, remained the capital of the new satrapy and also became the administrative headquarters for the satrapies of Assyria and Syria. One result was that the palace had to be enlarged.

Babylonia remained a wealthy and prosperous land, in contrast to Assyria, which was still a poor country. At the same time, the administration of the kingdom was more and more in the hands of the Persians, and the tax burdens grew heavier. This produced discontent, centring especially on the large temples in Babylon. Xerxes (486–465) had his residence in Babylon while he was crown prince, and he knew the country very well. When he assumed his kingship, he immediately curtailed the autonomy of the satrapies. This, in turn, gave rise to many rebellions. In Babylonia there were two short interim governments of Babylonian pretenders during 484–482. Xerxes retaliated by desecrating and partially destroying the holy places of the god Marduk and the Tower of Babel in the city of Babylon. Priests were executed, and the statue of Marduk was melted down.

The members of the royal family still resided in the palaces of the city of Babylon, but Aramaic became more and more the language of the official administration. One source of information for this period are the clay-tablet archives of the commercial house of Murashu and Sons of Nippur for the years of 455–403, which tell much about the important role the Iranians played in the country. The state domains were largely in their hands. They controlled many minor feudal tenants, grouped into social classes according to ancestry and occupation. The business people were predominantly Babylonians and Aramaeans, but there were also Jews.

The documents become increasingly sparse after 400. The cultural life of Babylon became concentrated in a few central cities, particularly Babylon and Uruk; Ur and Nippur were also important centres. The work of astronomers continued, as evidenced in records of observations. Naburimanni, living and working around 500, and Kidinnu, 5th or 4th century BC, were known to the Greeks; both astronomers are famous for their methods of calculating the courses of the Moon and the planets. In the field of literature, religious poetic works as well as texts of omens and Sumero-Akkadian word lists were constantly copied, often with commentaries. (W.T.v.S.)

Mesopotamia from c. 320 BC to c. AD 620

The political history of Mesopotamia between about 320 BC and AD 620 is divided among three periods of foreign rule—the Seleucids to 141 BC, the Parthians to AD 224,



Mesopotamia in Seleucid-Parthian times.

and the Sāsānians until the Arab invasions of the 7th century AD. Sources are scarce, consisting mainly of a few notices in the works of classical authors such as Strabo, Pliny, Polybius, and Ptolemy, while the cuneiform sources are mainly incantations, accounts of religious rites, and copies of ancient religious texts.

THE SELEUCID PERIOD

At the end of the Achaemenian Empire, Mesopotamia was partitioned into the satrapy of Babylonia in the south, while the northern part of Mesopotamia was joined with Syria in another satrapy. It is not known how long this division lasted, but, by the death of Alexander the Great in 323 BC, the north was removed from Syria and made a separate satrapy.

In the wars between the successors of Alexander, Mesopotamia suffered much from the passage and the pillaging of armies. When Alexander's empire was divided in 321 BC, one of his generals, Seleucus (later Seleucus I Nicator), received the satrapy of Babylonia to rule. From about 315 to about 312 BC, however, Antigonos I Monophthalmus (The "One-Eyed") took over the satrapy as ruler of all Mesopotamia, and Seleucus had to flee and accept refuge with Ptolemy of Egypt. With the aid of Ptolemy, Seleucus was able to enter Babylon in 312 BC (311 by the Babylonian reckoning) and hold it for a short time against the forces of Antigonos before marching to the east, where he consolidated his power. It is uncertain when he returned to Babylonia and reestablished his rule there; it may have been in 308, but by 305 BC he had assumed the title of king. With the defeat and death of Antigonos at the Battle of Ipsus in 301, Seleucus became the ruler of a large empire stretching from modern Afghanistan to the Mediterranean Sea. He founded a number of cities, the most important of which were Seleucia, on the Tigris, and Antioch, on the Orontes River in Syria. The latter, named after his father or his son, both of whom were called Antiochus, became the principal capital, while Seleucia became the capital of the eastern provinces. The dates of the founding of these

two cities are unknown, but presumably Seleucus founded Seleucia after he became king, while Antioch was built after the defeat of Antigonos.

Mesopotamia is scarcely mentioned in the Greek sources relating to the Seleucids, because the Seleucid rulers were occupied with Greece and Anatolia and with wars with the Ptolemies of Egypt in Palestine and Syria. Even the political division of Mesopotamia is uncertain, especially since Alexander, Seleucus, and Seleucus' son Antiochus I Soter all founded cities that were autonomous, like the Greek polis. The political division of the land into 19 or 20 small satrapies, which is found later, under the Parthians, began under the Seleucids. Geographically, however, Mesopotamia can be divided into four areas: Characene, also called Mesene, in the south; Babylonia, later called Asūristān, in the middle; northern Mesopotamia, where there was later a series of small states such as Gordyene, Osroene, Adiabene, and Garamaea; and finally the desert areas of the upper Euphrates, in Sāsānian times called Arabistān. These four areas had different histories down to the Arab conquest in the 7th century, although all of them were subject first to the Seleucids and then to the Parthians and Sāsānians. At times, however, several of the areas were fully independent, in theory as well as in fact, while the relations of certain cities with provincial governments and with the central government varied. From cuneiform sources it is known that traditional religious practices and forms of government as well as other customs continued in Mesopotamia; there were only a few Greek centres, such as Seleucia and the island of Ikaros (modern Faylakah, near Kuwait), where the practices of the Greek polis held sway. Otherwise, native cities had a few Greek officials or garrisons but continued to function as they had in the past.

Seleucia on the Tigris was not only the eastern capital but also an autonomous city ruled by an elected senate, and it replaced Babylon as the administrative and commercial centre of the old province of Babylonia. In the south several cities, such as Furat and Charax, grew rich

Seleucus

Geographic divisions

on the maritime trade with India; Charax became the main entrepôt for trade after the fall of the Seleucids. In the north there was no principal city, but several towns, such as Arbela (modern Irbil) and Nisibis (modern Nusaybin), later became important centres. In the desert region, "caravan cities" such as Hatra and Palmyra began their rise in the Seleucid period and had their heyday under the Parthians.

The only time that the Seleucid kings lost control of Mesopotamia was from 222 to 220 BC, when Molon, the governor of Media, revolted and marched to the west. When the new Seleucid king, Antiochus III, moved against him from Syria, however, Molon's forces deserted him, and the revolt ended. The Parthians, under their able king Mithradates I, conquered Seleucid territory in Iran and entered Seleucia in 141 BC. After the death of Mithradates I in 138 BC, Antiochus VII began a campaign to recover the Seleucid domains in the east. This campaign was successful until Antiochus VII lost his life in Iran in 129 BC. His death ended Seleucid rule in Mesopotamia and marked the beginning of small principalities in both the south and north of Mesopotamia.

Seleucid rule brought changes to Mesopotamia, especially in the cities where Greeks and Macedonians were settled. In these cities the king usually made separate agreements with the Greek officials of the city regarding civil and military authority, immunity from taxes or *corvée*, or the like. Native cities continued with their old systems of local government, much as they had under the Achaemenians. Greek gods were worshiped in temples dedicated to them in the Greek cities, and native Mesopotamian gods had temples dedicated to them in the native cities. In time, however, syncretism and identification of the foreign and local deities developed. Although the policy of Hellenization was not enforced upon the population, Greek ideas did influence the local educated classes, just as local practices were gradually adopted by the Greeks. As in Greece and the lands of the eastern Mediterranean, in Mesopotamia the philosophies of the Stoics and other schools probably had an impact, as did mystery religions; both were hallmarks of the Hellenistic Age. Unfortunately there is no evidence from the east on the popularity of Greek beliefs among the local population, and scholars can only speculate on the basis of the fragmentary notices in authors such as Strabo. The Seleucid rulers respected the native priesthoods of Mesopotamia, and there is no record of any persecutions. On the contrary, the rulers seem to have favoured local religious practices, and ancient forms of worship continued. Cuneiform writing by priests, who copied incantations and old religious texts, continued into the Parthian period.

The administrative institutions of the countryside of Mesopotamia remained even more traditional than those of the cities; the old taxes were simply paid to new masters. The satrapy, much reduced in size from Achaemenian times, was the basis for Seleucid control of the countryside. A satrap or strategus (a military title) headed each satrapy, and the satrapies were divided into hyparchies or eparchies; the sources that use these and other words, such as *toparchy*, are unclear about the subdivisions of the satrapy. There was a great variety of smaller units of administration. In the capital and in the provincial centres, both Greek and Aramaic were used as the written languages of the government. The use of cuneiform in government documents ceased sometime during the Achaemenian period, but it continued in religious texts until the 1st century of the Common era. The archives were managed both in the capital and in provincial cities by an official called a *bibliophylax*. There were many financial officials (*oikonomoi*); some of them oversaw royal possessions, and others managed local taxes and other economic matters. The legal system in the Seleucid empire is not well understood, but presumably both local Mesopotamian laws and Greek laws, which had absorbed or replaced old Achaemenian imperial laws, were in force. Excavations at Seleucia have uncovered thousands of seal impressions on clay, evidence of a developed system of controls and taxes on commodities of trade. Many of the sealings are records of payment of a salt tax. Most of the

tolls and tariffs, however, were local assessments rather than royal taxes.

Artistic remains from the Seleucid period are exceedingly scarce, and, in contrast to Achaemenian art, no royal or monumental art has been recovered. One might characterize the objects that can be dated to the Seleucid era as popular or private art, such as seals, statuettes, and clay figurines. Both Greek and local styles are found, with an amalgam of styles prevalent at the end of Seleucid rule, evidence of a syncretism in cultures. The numerous statues and statuettes of Heracles found in the east testify to the great popularity of the Greek deity, in Mesopotamia identified with the local god Nergal.

Aramaic was the "official" written language of the Achaemenian Empire; after the conquests of Alexander the Great, Greek, the language of the conquerors, replaced Aramaic. Under the Seleucids, however, both Greek and Aramaic were used throughout the empire, although Greek was the principal language of government. Gradually Aramaic underwent changes in different parts of the empire, and in Mesopotamia in the time of the Parthians it evolved into Syriac, with dialectical differences from western Syriac, used in Syria and Palestine. In southern Mesopotamia, other dialects evolved, one of which was Mandaic, the scriptural language of the Mandaean religion.

Literature in local languages is nonexistent, except for copies of ancient religious texts in cuneiform writing and fragments of Aramaic writing. There were authors who wrote in Greek, but little of their work has survived and that only as excerpts in later works. The most important of these authors was Berosus, a Babylonian priest who wrote about the history of his country, probably under Antiochus I (reigned 281–261 BC). Although the excerpts of his work that are preserved deal with the ancient, mythological past and with astrology and astronomy, the fact that they are in Greek is indicative of interest among local Greek colonists in the culture of their neighbours. Another popular author was Apollodorus of Artemita (a town near Seleucia), who wrote under the Parthians a history of Parthia in Greek as well as other works on geography. Greek continued to be a lingua franca used by educated people in Mesopotamia well into the Parthian period.

Under the Seleucid system of dating, as far as is known, a fixed year became the basis for continuous dating for the first time in the Middle East. The year chosen was the year of entry of Seleucus into Babylon, 311 BC according to the Mesopotamian reckoning and 312 BC according to the Syrians. Before this time, dating had been only according to the regnal years of the ruling monarch (*e.g.*, "fourth year of Darius"). The Parthians, following the Seleucids, sought to institute their own system of reckoning based on some event in their past that scholars can only surmise—possibly the assumption of the title of king by the first ruler of the Parthians, Arsaces.

Since Greece was overpopulated at the beginning of Seleucid rule, it was not difficult to persuade colonists to come to the east, especially when they were given plots of land (*cleroii*) from royal domains that they could pass on to their descendants; if they had no descendants, the land would revert to the king. Theoretically all land belonged to the ruler, but actually local interests prevailed. As time passed, however, the influx of Greek colonists diminished and then ended when the wars of the Hellenistic kings interrupted this movement. Nonetheless, Greek influences continued, and it is fascinating to find in cuneiform documents records of families where the father has a local name and his son a Greek one, and vice versa. Inasmuch as Mesopotamia was peaceful under the Seleucids, the processes of accommodation and assimilation among the people appear to have flourished.

THE PARTHIAN PERIOD

The coming of the Parthians changed Mesopotamia even less than the establishment of the Seleucid kingdom had, for as early as the middle of the 2nd century BC local dynasts had proclaimed their independence. There is no evidence indicating whether the cities of Mesopotamia surrendered piecemeal or all at once or whether they submitted voluntarily or after fighting. In any case, Seleu-

The influence of Greek philosophy and religion

Literature

cia was treated better by the Parthians than it had been by the Seleucids, and the local government retained its autonomy. Parthian troops did not occupy Seleucia but remained in a garrison site called Ctesiphon near Seleucia; it later grew into a city and replaced Seleucia as the capital. In Characene in southern Mesopotamia a Seleucid satrap with an Iranian name, Hyspaosines, issued coins about 125 bc, a sign of his independence; the actual date for this may have been earlier. He changed the name of the city Antiochia on the lower Tigris to Spasinou Charax, meaning "The Fort of Hyspaosines," and made it his capital. All the coins issued from his capital have Greek legends. His troops moved north and occupied Babylon and Seleucia probably sometime in 127 bc, when the Parthians were fighting nomadic invaders in the eastern part of their territory. His rule there must have been short, however, for the Parthian governor of Babylon and the north, Himerus, was back in Seleucia and Babylon by 126. Himerus could not have been a rebel, since he struck coins in the name of the Parthian rulers Phraates II and Artabanus II, both of whom were killed in fighting in eastern Iran. Himerus abused his power and is said to have oppressed the cities of Mesopotamia, plundering them and killing their inhabitants. Cuneiform documents from Babylon stop after this date, indicating that the city did not survive the depredations of Himerus. He vanished, however, and Parthian sovereignty was restored by the ninth Arsacid king, Mithradates II, who came to the throne about 124 bc; he was the son of Artabanus II. Mithradates II recovered all Mesopotamia and conquered Characene, overstriking coins of Hyspaosines and driving him from his capital in 122 or 121 bc. By 113, if not earlier, Dura-Europos on the Euphrates was in Parthian hands. In 95 bc the Armenian Tigranes II, a hostage at the court of Mithradates, was placed on the throne of Armenia by his Parthian overlord, and the small kingdoms of northern Mesopotamia—Adiabene, Gordyene, and Osroene—gave allegiance to Mithradates. Mithradates II died about 87 bc, although he may have died earlier, since the period after 90 bc is dark and a usurper named Gotarzes may have ruled for a few years in Mesopotamia. During the reign of Mithradates II the first contacts with Rome, under Lucius Cornelius Sulla, were made, and portents of future struggles were evident in the lack of any agreement between the two powers. Sulla was sent to the east by the Roman Senate to govern Cilicia in Anatolia. In 92 bc Orobazes, an ambassador from Mithradates II, came to him seeking a treaty, but nothing was concluded, since instructions from Rome did not include negotiations with the Parthian power.

Tigranes II took advantage of struggles between several claimants to the Parthian throne to expand Armenian territory into Mesopotamia, and the small states in the north gave him their allegiance. It was not until 69 bc, when the Roman general Lucius Licinius Lucullus captured Tigranokerta, Tigranes' capital, that Mesopotamia returned to Parthian rule. Thereafter wars between the Romans and the Parthians were to dominate the political history of Mesopotamia.

The Parthians left the local administrations and rulers intact when they conquered Mesopotamia. According to Pliny the Elder (*Natural History* VI. 112) the Parthian empire consisted of 18 kingdoms, 11 of which were called the upper kingdoms (or satrapies), while 7 were called lower kingdoms, meaning that they were located on the plains of Mesopotamia. The centre of the lower kingdoms was ancient Babylonia, called Beth Aramaye in Aramaic, and it was governed directly by the Parthian ruler. In the south was Characene, while to the northeast of Ctesiphon, which had supplanted Seleucia as the Parthian capital, was Garama, with its capital at modern Kirkūk. Adiabene had Arbela as its capital, and farther north was a province called Beth Nuhadra in Aramaic, which seems to have been governed by a general who was directly responsible to the Parthian king, because this province bore the brunt of Roman invasions. Nisibis was the main city of the desert area of Arabistān, but at the end of the Parthian period the desert caravan city of Hatra claimed hegemony over this area. There were other principalities

in the northwest: Sophene, where Tigranes' capital was located; Gordyene and Zabdicene (near modern Çölemerik in eastern Turkey), located to the east of Sophene; and Osroene, with its capital Edessa (modern Urfa, Tur.), which lay inside the Roman sphere of influence. Rule over so many small kingdoms gave Mithradates II the title "King of Kings," also borne by later Parthian rulers.

The defeat of the Roman legions under Marcus Licinius Crassus by the Parthians at the Battle of Carrhae (Carrhae is the Roman name for Harran) in 53 bc heralded a period of Parthian power and expansion in the Middle East, but the tide turned under Mark Antony in 36–34 bc, and thereafter the power structure in the east remained volatile, with the two great states, Rome and Parthia, contending for predominance in the region. Armenia was a perennial bone of contention between the two powers, each of which sought to put its candidate on the throne.

Parthian rule was not firm over all Mesopotamia; thus, for example, during the reign of Artabanus III (AD 12–38) the Jewish brigands Asinaeus and Anilaeus set up a free state north of Ctesiphon that lasted 15 years before it was overcome by the Parthians. With the end of cuneiform records and with the attention of classical sources turned to the wars between the Romans and the Parthians, information about internal affairs in Mesopotamia becomes almost nonexistent. Hellenism continued to flourish, for many Parthian kings had the epithet "Philhellene" placed on their coins, but during the last two centuries of Parthian rule Greek influences declined in favour of Iranian ones, while central authority suffered from the usurpations of powerful nobles and local kings. From coinage it is known that the city of Seleucia revolted against central control at the end of Artabanus' reign and maintained its independence for a number of years. Peace was broken by the Roman emperor Nero, who sought to put his client on the throne of Armenia, but, after several years of conflict, peace was arranged in 63. Vologeses I (c. AD 51–80) founded the city Vologesias, near Seleucia, as his capital, but the whole area (including Ctesiphon and Seleucia) became an urban complex called Māhōzē in Aramaic and Al-Madā'in in Arabic; both names mean "The Cities."

Internal rivalries in the Parthian state gave the Romans an opportunity to attack, and control over Armenia was the *casus belli* for the Roman emperor Trajan's advance into Mesopotamia in 116. Adiabene, as well as the entire Tigris-Euphrates basin of northern Mesopotamia, was incorporated as a province into the Roman Empire. Trajan advanced to the Persian Gulf, but he died of illness and his successor Hadrian made peace, abandoning the conquests in Mesopotamia, although client states remained.

The second century of the Common era was a dark period in Parthian history, but it was a time of growth in wealth and influence of the caravan cities of Palmyra, Hatra, and Mesene (formerly Characene, situated at the confluence of the Tigris and Euphrates). Armenia continued to be a bone of contention between the two great powers, and hostilities occasionally flared up. In 164–165 the Roman general Gaius Avidius Cassius captured the capital cities Ctesiphon and Seleucia, but an epidemic forced the Romans to retreat and peace was restored. Returning soldiers spread the disease throughout the Roman Empire, with devastating consequences. The terms of peace favoured the Romans, who secured control of Nisibis and the Khābūr River valley. The next great war was the invasion of the Roman emperor Septimius Severus to punish the Parthians, who had supported his rival Pescennius Niger and had annexed some territory in Mesopotamia in return for their support. Severus took and sacked Ctesiphon in 198. Because the devastated countryside contained no supplies for the Romans, they were soon compelled to retreat. A siege of Hatra in 199 by Severus failed, and peace was made. Conflict between two claimants to the Parthian throne, Vologeses IV or V and Artabanus V, gave the Roman emperor Caracalla an excuse to invade Adiabene, but in 217 he was assassinated on the road from Edessa to Carrhae, and the Romans made peace. The end of the Parthian kingdom was near, and the advent of the Sāsānians brought a new phase in the history of Mesopotamia.

Parthian rule brought little change in the administration

Mithradates II

Conflict between Rome and Parthia

and institutions of Mesopotamia as it had existed under the Seleucids, except for a general weakening of central authority under the feudal Parthians. The Parthians instituted a new era, beginning in 247 BC, but it paralleled rather than replaced the Seleucid era of reckoning, and the Parthian vanished at the end of the dynasty. As far as can be determined, Hellenism was never proscribed under the Parthians, although it grew weaker toward the end of Parthian rule. From archaeological surveys around Susa, located in the kingdom of Elymais in modern Khūzestān, and from the Diyālā plain northeast of Ctesiphon, it seems that the population of the land increased greatly under the Parthians, as did trade and commerce. The coinage of the later Parthian rulers became more and more debased, probably as a result of the many internecine wars and the lack of control by the central authority. Local rulers also issued their own coinages in Persis, Elymais, Mesene, and elsewhere.

Demographic changes

Changes took place in the demography of Mesopotamia under the Parthians, and perhaps the most striking development among the population was the increase of Arab infiltration from the desert, which resulted in Arab dynasties in the oasis settlements of Palmyra and Hatra. Similarly, an influx of Armenian settlers in the north changed the composition of the local population. After the fall of the Temple of Jerusalem to the Romans in 70, many Jews fled to Mesopotamia, where they joined their coreligionists; Nehardea, north of Ctesiphon, became a centre of Jewish population. Naturally also many migrants from the east came to Mesopotamia in the wake of the Parthian occupation. With many merchants from east and west passing through or remaining in Mesopotamia, the population became more diverse than it had previously been.

During the Parthian occupation the ancient religion and cults of Mesopotamia came to an end and were replaced by mixed Hellenic and Oriental mystery religions and Iranian cults. Local Semitic cults of Bel, Allat, and other deities flourished alongside temples dedicated to Greek gods such as Apollo. The sun deity Shamash was worshipped at Hatra and elsewhere, but the henotheism of the ancient Middle East was giving way to acceptance of universalist religions, if the prevalent view cannot yet be called one of monotheism. In Mesopotamia, in particular, the influence of Jewish monotheism, with the beginning of rabbinic schools and the organization of the community under a leader, the exilarch (*resh galuta* in Aramaic), must have had a significant influence on the local population. Toward the end of the reign of Artabanus III, the royal family of Adiabene converted to Judaism. In the first two centuries of the Common era, Christianity and various baptismal sects also began to expand into Mesopotamia. So far no Mithraeums (underground temples for the worship of the god Mithra), such as existed in the Roman Empire, have been found in Mesopotamia, except at Dura-Europus, where Roman troops were stationed. Many local cults and shrines, such as that of the Sabians and their moon deity at Harran, however, continued to exist until the Islāmic conquest. Parthian Zoroastrianism reinforced local Zoroastrian communities in Mesopotamia left from the time of the Achaemenians, and one of the Gnostic baptismal religions, Mandaeanism, which is still in existence, had its beginning at this time. Although Christian missionaries were active in Mesopotamia in the Parthian period, no centres, such as the one established later at Nisibis, have been reported, and it may be supposed that their activity at first was mainly confined to Jewish communities.

Art and architecture

Archaeological evidence indicates that the Parthians had a more marked influence on art and architecture. Local schools of art flourished, and at first Greek ideals predominated, but in the last two centuries of Parthian rule a "Parthian style" is evident in the art recovered from Mesopotamia and other regions. Whereas Achaemenian and Sāsānian art are royal or imperial and monumental, Parthian art, like Seleucid art, can be characterized as "popular." Parthian works of art reflect the many currents of culture among the populace, and one may say that it is expressionist and stylized, in contrast with Greek and Roman naturalistic or realistic art. The characteristics of

Parthian art in Mesopotamia are total frontality (*i.e.*, the representation of figures in full face) in portraits, along with an otherworldly quality. In Middle Eastern art from previous periods, figures were almost always shown in profile. Another new feature of Parthian art is the frequent portrayal of the "flying gallop" in sculpture and painting, not unexpected in view of the importance of cavalry and mounted archers in the Parthian armies. Likewise, Parthian costume, with baggy trousers, became the mode over much of the Middle East and is portrayed in painting and sculpture. In architecture the use of *ayvans* (arches in porticoes) and domed vaults is attributed to the Parthian period; they may have originated in Mesopotamia. Parthian art influenced that of the Nabataeans in Roman territory, as it did others throughout the Middle East.

Parthian was an Iranian language written in the Aramaic alphabet. It had an enormous number of words and even phrases that were borrowed from Aramaic, and scribal training was necessary to learn these. Syriac, being a Semitic language with emphasis on consonants, evolved several alphabets based on the Aramaic alphabet. The Aramaic alphabet was better suited to Syriac than to Parthian phonology. Parthian was therefore difficult to read and was mainly used by scribes or priests for official or religious writings.

The largest lacuna is in literature from the Parthian period. The largely oral literature of the Parthians, famous for their minstrels and poetry, does not seem to have found many echoes in Mesopotamia, where the settled society contrasted with the heroic, chivalric, and feudal society of the Iranian nomads that continued to dominate Parthian mores even after they had settled in Mesopotamia. Nonetheless, the end of the Parthian period saw the beginning of Syriac literature, which is Christian Aramaic, and some of early Syriac literature, such as the "Song of the Pearl," contains Parthian elements. In the realm of language, rather than literature, the writing of Aramaic changes to Parthian in the 2nd century AD, as can be seen from a bilingual (Greek and Parthian) inscription on a bronze statue from Seleucia dated AD 150–151. It tells how Vologeses III defeated the king of Mesene and took over the entire country. After this period one no longer speaks of Aramaic, but of Parthian and Syriac written in a new cursive alphabet.

THE SĀSĀNIAN PERIOD

The Sāsānian period marks the end of the ancient and the beginning of the medieval era in the history of the Middle East. Universalist religions such as Christianity, Manichaeism, and even Zoroastrianism and Judaism absorbed local religions and cults at the beginning of the 3rd century. Both the Sāsānian and the Roman empires ended by adopting an official state religion, Zoroastrianism for the former and Christianity for the latter. In Mesopotamia, however, older cults such as that of the Mandaean, the moon cult of Harran, and others continued alongside the great religions. The new rulers were not as tolerant as the Seleucids and Parthians had been, and persecutions occurred under Sāsānian rule.

After Ardashir I, the first of the Sāsānians, consolidated his position in Persis (modern Fārs province), he moved into southern Mesopotamia, and Mesene submitted. In 224 he defeated and killed the last Parthian ruler, Artabanus V, after which Mesopotamia quickly fell before him and Ctesiphon became the main capital of the Sāsānian empire. In 230 Ardashir besieged Hatra but failed to take it. Hatra called on Roman aid, and in 232 the Roman emperor Severus Alexander launched a campaign that halted Ardashir's progress. At the death of Severus Alexander in 235 the Sāsānians took the offensive, and probably in 238 Nisibis and Harran came under their control. Hatra was probably captured in early 240, after which Ardashir's son Shāpūr was made coregent; Ardashir himself died soon afterward. The Roman emperor Gordian III led a large army against Shāpūr I in 243. The Romans retook Harran and Nisibis and defeated the Sāsānians at a battle near Resaina, but at Anbār, renamed Pērōz-Shāpūr ("Victorious Is Shāpūr"), the Sāsānians inflicted a defeat on the Romans, who lost their emperor. His successor, Philip the Arabian,

Ardashir I

made peace, giving up Roman conquests in northern Mesopotamia. Osroene, however, which had been returned to the local ruling family of Abgar by Gordian, remained a vassal state of the Romans. Shāpūr renewed his attacks and took many towns, including Dura-Europus, in 256 and later moved into northern Syria and Anatolia. The defeat and capture of the Roman emperor Valerian at the gates of Edessa, probably in 259, was the high point of his conquests in the west. On Shāpūr's return to Ctesiphon the ruler of Palmyra, Septimius Odaenathus (also called Odaenath), attacked and defeated his army, seizing booty. Odaenathus took the title of emperor, conquered Harran and Nisibis, and threatened Ctesiphon in 264–266. His murder relieved the Sāsānians, and in 273 the Roman emperor Aurelian sacked Palmyra and restored Roman authority in northern Mesopotamia. Peace between the two empires lasted until 283, when the Roman emperor Carus invaded Mesopotamia and advanced on Ctesiphon, but the Roman army was forced to withdraw after Carus' sudden death. In 296 Narseh I, the seventh Sāsānian king, took the field and defeated a Roman force near Harran, but in the following year he was defeated and his family was taken captive. As a result, the Romans secured Nisibis and made it their strongest fortress against the Sāsānians. The Roman province of Mesopotamia, which was the land between the Euphrates and Tigris in the northern foothills, became in effect a military area with *limes* (the fortified frontiers of the Roman Empire) and highly fortified towns.

Under Shāpūr II the Sāsānians again took the offensive, and the first war lasted from 337 to 350; it ended with no result as Nisibis was successfully defended by the Romans. In 359 Shāpūr again invaded Roman territory and captured the Roman fortress Amida after a long and costly siege. In 363 the emperor Julian advanced almost to Ctesiphon, where he died, and his successor Jovian had to give up Nisibis and other territories in the north to the Sāsānians. The next war lasted from 502 to 506 and ended with no change. War broke out again in 527, lasting until 531, and even the Byzantine general Belisarius was not able to prevail; as usual, the boundaries remained unchanged. In 540 the Sāsānian king Khosrow (Khosroes) I invaded Syria and even took Antioch, although many fortresses behind him in northern Mesopotamia remained in Byzantine hands. After much back-and-forth fighting, peace was made in 562. War with the Byzantine Empire resumed 10 years later, and it continued under Khosrow's successor, Hormizd IV. Only in 591, in return for their assistance in the restoration to the Sāsānian throne of Khosrow II, who had been deposed and had fled to Byzantine territory, did the Byzantines regain territory in northern Mesopotamia. With the murder in 602 of the Byzantine emperor Maurice, who had been Khosrow's benefactor, and the usurpation of Phocas, Khosrow II saw a golden opportunity to enlarge Sāsānian domains and to take revenge for Maurice. Persian armies took all northern Mesopotamia, Syria, Palestine, Egypt, and Anatolia. By 615, Sāsānian forces were in Chalcedon, opposite Constantinople. The situation changed completely with the new Byzantine emperor Heraclius, who, in a daring expedition into the heart of enemy territory in 623–624, defeated the Sāsānians in Media. In 627–628 he advanced toward Ctesiphon, but, after sacking the royal palaces at Dastagird, northeast of Ctesiphon, he retreated.

After the death of Khosrow II, Mesopotamia was devastated not only by the fighting but also by the flooding of the Tigris and Euphrates, by a widespread plague, and by the swift succession of Sāsānian rulers, which caused chaos. Finally in 632 order was restored by the last king, Yazdegerd III, but in the following year the expansion of the Muslim Arabs began and the end of the Sāsānian empire followed a few years afterward.

Unlike the Parthians, the Sāsānians established their own princes as rulers of the small kingdoms they conquered, except on the frontiers, where they accepted vassals or allies because their hold over the frontier regions was insecure. By placing Sāsānian princes over the various parts of the empire, the Sāsānians maintained more control than the Parthians had. The provincial divisions were more

systematized, and there was a hierarchy of four units—the satrapy (*shahr* in Middle Persian), under which came the province (*ōstan*), then a district (*tassug*), and finally the village (*deh*). In Mesopotamia these divisions were changed throughout Sāsānian history, frequently because of Roman invasions.

Many native tax collectors were replaced by Persians, who were more trusted by the rulers. In addition to the many tolls and tariffs, corvée, and the like, the two basic taxes were the land and poll taxes. The latter were not paid by the nobility, soldiers, civil servants, and the priests of the Zoroastrian religion. The land tax was a percentage of the harvest, but it was determined before the collection of the crops, which naturally caused many problems. Khosrow I undertook a new survey of the land and imposed the tax in a prearranged sum based on the amount of cultivable land, the quantity of date palms and olive trees, and the number of people working on the land. Taxes were to be paid three times a year. Abuses were still rampant, but this was better than the old system; at least, if a drought or some other calamity occurred, taxes could be reduced or remitted. Although information is contradictory, it appears that religious communities other than the Zoroastrian one had extra taxes imposed on them from time to time. This was especially true of the growing Christian community, particularly in the time of Shāpūr II, after Christianity became the official religion of the Roman Empire.

Religious communities became fixed under the Sāsānians, and Mesopotamia with its large Jewish and Christian populations experienced changes because of the shift in primary allegiance from the ruler to the head of the religious group. The exilarch of the Jews had legal and tax-collecting authority over the Jews of the Sāsānian empire. Mani, the founder of the Manichaean religion, was born in lower Mesopotamia, and his religion spread quickly both to the east and west, even before his death. In its homeland, Mesopotamia, it came under severe persecution by the priests of the Zoroastrian religion, who viewed Manichaeism as a dangerous heresy. Christianity, however, was viewed not as a heresy but as a separate religion, tolerated until it became the official religion of the enemy Roman Empire; Christians were then regarded as potential traitors to the Sāsānian state. The first large growth of Christianity in Mesopotamia came with the deportation and resettlement of Christians, especially from Antioch with its patriarch, during Shāpūr I's wars with the Romans. In a synod convened in 325, the metropolitan see of Ctesiphon was made supreme over other sees in the Sāsānian empire, and the first patriarch or catholicos was Papa. In 344 the first persecutions of Christians began; they lasted with varying degrees of severity until 422, when a treaty with the government ended the persecutions.

The earliest contemporary mention of Christians in Mesopotamia is in the inscriptions of Kartēr, the chief Zoroastrian priest after the reign of Shāpūr I. He mentions both Christians and Nazareans, possibly two kinds of Christians, Greek-speaking and Syriac-speaking, or two sects. It is not known which groups are meant, but it is known that followers of the Gnostic Christian leaders Bardesanes (Bar Daišān) and Marcion were active in Mesopotamia. Later, after the Nestorian church separated from the Monophysites, whose centre was in Antioch, the Nestorian church dominated Mesopotamia until the end of the Sāsānian dynasty, when the Monophysites were growing in numbers. After about 485 the Sāsānian government was satisfied that the Nestorian church in their domains was not loyal to Byzantium, and further persecutions were not state-inspired but rather prosecuted by the Zoroastrian clergy. At the end of the Sāsānian period, the Nestorians were fighting the Monophysites, now called Jacobites, more than the Zoroastrians. The Jacobites established many monasteries, especially in northern Mesopotamia, whereas the Nestorians were cool toward monasticism.

Ethnicity became less important than religious affiliation under the Sāsānians, who thus changed the social structure of Mesopotamia. The Arabs continued to grow in numbers, both as nomads and as settled folk, and Ara-

Taxation

Wars
with the
Byzantine
EmpireGrowth of
the Arab
population

bic became widely spoken. King Nu'mān III of the Arab client kingdom of the Lakhmids of Al-Hīrah in southern Mesopotamia became a Christian in 580, but in 602 he was deposed by Khosrow II, who made the kingdom a province of the empire. This act removed a barrier against inroads by Arab tribesmen from the desert, and, after the union of Arabs in the peninsula under the banner of Islām, the fate of the Sāsānian empire was sealed. The Muslims, on the whole, were welcomed in Mesopotamia as deliverers from the foreign yoke of the Persians, but the conversion of the mass of the population to Islām did not proceed rapidly, mainly because of the well-organized Christian and Jewish communities. The arrival of Islām, of course, changed the history of Mesopotamia more than any other event in its history. (R.N.F.)

BIBLIOGRAPHY

General works. *The Cambridge Ancient History* contains much relevant information, especially vol. 1–2, 3rd ed. (1970–75), vol. 3–4, 2nd ed. (1982–88), and vol. 6 (1927); they include lengthy and richly documented chapters covering Mesopotamian prehistory to the time of Alexander the Great's conquest of the region. Chapters on Mesopotamia under the Seleucids, Parthians, and Sāsānians are included in *The Cambridge History of Iran*, vol. 3 (1983). ROBERT MCC. ADAMS, *The Land Behind Baghdad: A History of Settlement on the Diyala Plains* (1965); NICHOLAS POSTGATE, *The First Empires* (1977); GEORGES ROUX, *Ancient Iraq*, 2nd ed. (1980); RICHARD N. FRYE, *The History of Ancient Iran* (1984); SETON LLOYD, *The Archaeology of Mesopotamia: From the Old Stone Age to the Persian Conquest*, rev. ed. (1984); and MICHAEL ROAF, *Cultural Atlas of Mesopotamia and the Ancient Near East* (1990), also provide broad coverage. I.M. DIAKONOV (ed.), *Ancient Mesopotamia: Socio-Economic History*, trans. from Russian (1969, reissued 1981), collects representative articles by Diakonov and others on Mesopotamian history, with emphasis on social and economic aspects. A. LEO OPPENHEIM, *Ancient Mesopotamia: Portrait of a Dead Civilization*, rev. ed. completed by ERICA REINER (1977), includes some controversial views.

Prehistory to the Old Babylonian period. A balanced picture of political, social, and economic history may be found in JEAN BOTTÉRO, ELENA CASSIN, and JEAN VERCOUTTER (eds.), *The Near East: The Early Civilizations* (1968; originally published in German, 3 vol., 1965–67), with contributions on prehistory and protohistory, Akkad, Early Dynastic history, the 3rd dynasty of Ur, and the Old Babylonian period. ADAM FALKENSTEIN, *The Sumerian Temple City*, trans. from French (1974), is a very short work describing the Sumerian temple economy and its political implications. DIETZ OTTO EDZARD, *Die zweite Zwischenzeit Babyloniens* (1957), offers details on the history of the Old Babylonian period from the 3rd dynasty

of Ur to the end of Hammurabi. MOGEN TROLLE LARSEN, *The Old Assyrian City-State and Its Colonies* (1976), is a standard work on the Old Assyrian trade colonies in Anatolia. GERNOT WILHELM, *The Hurrians* (1989; originally published in German, 1982), is the best book on the third cultural element in early Mesopotamian history. FIORELLA IMPARATI, *I Hurriti* (1964), offers a short synopsis. HANS J. NISSEN, *Mesopotamia Before 5000 Years* (1987), includes a comprehensive bibliography for the early periods.

Mesopotamia to the end of the Achaemenian period. Histories of Assyria and Babylonia include WOLFRAM VON SODEN, *Einführung in die Altorientalistik* (1985); J.A. BRINKMAN, *Prelude to Empire: Babylonian Society and Politics, 747–626 B.C.* (1984); STEFAN ZAWADZKI, *The Fall of Assyria and Median-Babylonian Relations* (1988); and H.W.F. SAGGS, *The Might That Was Assyria* (1984, reprinted 1990), and *The Greatness That Was Babylon: A Survey of the Ancient Civilization of the Tigris-Euphrates Valley* (1988). JOAN OATES, *Babylon*, rev. ed. (1986), deals with history and civilization. WOLFRAM VON SODEN, *Herrscher im alten Orient* (1954), examines Assyrian and Babylonian politics. Standard works, now partly out-of-date, include A.T. OLMSTEAD, *History of Assyria* (1923, reprinted 1975); and BRUNO MEISSNER, *Babylonien und Assyrien*, 2 vol. (1920–25). J.A. BRINKMAN, *A Political History of Post-Kassite Babylonia, 1158–722 B.C.* (1968), is an extensive special study, with complete documentation. MUHAMMAD A. DANDAMAEV, *Slavery in Babylonia: From Nabopolassar to Alexander the Great (626–331 B.C.)*, rev. ed. (1984; originally published in Russian, 1974), includes an extensive bibliography. JACOB NEUSNER, *A History of the Jews in Babylonia*, 5 vol. (1965–70), studies in detail the history of the Jews in Mesopotamia.

Mesopotamia from c. 320 BC to c. AD 620. GETZEL M. COHEN, *The Seleucid Colonies: Studies in Founding, Administration, and Organization* (1978), discusses the relationship of the central government with provinces and client states. Details of Seleucid rule may be found in MAURICE MEULEAU, "Mesopotamia Under the Seleucids," in PIERRE GRIMAL (ed.), *Hellenism and the Rise of Rome* (1968; originally published in German, 1965), pp. 266–289; and in the essays in AMÉLIE KUHR and SUSAN SHERWIN-WHITE (eds.), *Hellenism in the East: Interaction of Greek and Non-Greek Civilizations after Alexander's Conquest* (1987). NELSON C. DEBEVOISE, *A Political History of Parthia* (1938, reissued 1968), is a standard political history of the Arsacid dynasty. SHELDON ARTHUR NODELMAN, "A Preliminary History of Characene," *Berytus*, 13(2):83–121 (1960), is the standard history of Characene. LOUIS DILLEMANN, *Haute Mésopotamie orientale et pays adjacents* (1962), offers a historical geography of northern Mesopotamia with many details. J.M. FIE, *Assyrie chrétienne*, 3 vol. (1965–68), provides a historical geography of the Christian communities of northern Mesopotamia from Syriac sources.

(D.O.E./W.T.v.S./R.N.F.)

Metabolism

Living organisms are unique in that they can extract energy from their environments and use it to carry out activities such as movement, growth and development, and reproduction. But how do living organisms—or, their cells—extract energy from their environments, and how do cells use this energy to synthesize and assemble the components from which the cells are made?

The answers to these questions lie in the enzyme-mediated chemical reactions that take place in living matter (metabolism). Hundreds of coordinated, multistep reac-

tions, fueled by energy obtained from nutrients and/or solar energy, ultimately convert readily available materials into the molecules required for growth and maintenance.

The physical and chemical properties of the components of living things dealt with in this article are found in the articles BIOCHEMICAL COMPONENTS OF ORGANISMS; CELLS; and PHOTOSYNTHESIS.

For coverage of related topics in the *Macropædia* and *Micropædia*, see the *Propædia*, section 322, and the *Index*.

The article is divided into the following sections:

-
- A summary of metabolism 893
 - The unity of life 893
 - Biological energy exchanges 894
 - The carrier of chemical energy
 - Catabolism
 - Anabolism
 - Integration of catabolism and anabolism 896
 - Fine control
 - Coarse control
 - The study of metabolic pathways 896
 - The fragmentation of complex molecules 897
 - The catabolism of glucose 898
 - Glycolysis
 - The phosphogluconate pathway
 - The catabolism of sugars other than glucose 900
 - The catabolism of lipids (fats) 901
 - Fate of glycerol
 - Fate of fatty acids
 - The catabolism of proteins 903
 - Removal of nitrogen
 - Disposal of nitrogen
 - Oxidation of the carbon skeleton
 - The combustion of food materials 905
 - The oxidation of molecular fragments 905
 - The oxidation of pyruvate
 - The tricarboxylic acid (TCA) cycle
 - Biological energy transduction 907
 - Adenosine triphosphate as the currency of energy exchange
 - Energy conservation
 - The biosynthesis of cell components 909
 - The nature of biosynthesis 909
 - The stages of biosynthesis
 - Utilization of ATP
 - The supply of biosynthetic precursors 910
 - Anaplerotic routes
 - Growth of microorganisms on TCA cycle intermediates
 - The synthesis of building blocks 911
 - Gluconeogenesis
 - Lipid components
 - Amino acids
 - Mononucleotides
 - The synthesis of macromolecules 916
 - Carbohydrates and lipids
 - Nucleic acids and proteins
 - Regulation of metabolism 919
 - Fine control 919
 - End-product inhibition
 - Positive modulation
 - Energy state of the cell
 - Coarse control 920
 - Metabolic diseases 921
 - General considerations 921
 - The derivation of specific metabolic disorders 921
 - Disorders of carbohydrate metabolism 922
 - Glycogen storage diseases
 - Galactosemia
 - Fructose disorders
 - Pyruvate disorders
 - Intestinal carbohydrate malabsorption
 - Mucopolysaccharidoses
 - Disorders of lipid metabolism 924
 - Blood lipid disorders
 - Lipid oxidative disorders
 - Tissue lipid disorders
 - Disorders of amino-acid metabolism 924
 - Phenylalanine and tyrosine
 - Urea cycle enzymes
 - Branched chain amino acids
 - Disorders of amino-acid transport
 - Disorders of porphyrin metabolism 927
 - Disorders of purine and pyrimidine metabolism 927
 - Bibliography 928
-

A summary of metabolism

THE UNITY OF LIFE

At the cellular level of organization, the main chemical processes of all living matter are similar, if not identical. This is true for animals, plants, fungi, or bacteria; where variations occur (such as, for example, in the secretion of antibodies by some molds), the variant processes are but variations on common themes. Thus, all living matter is made up of large molecules called proteins, which provide support and coordinated movement, as well as storage and transport of small molecules, and, as catalysts, enable chemical reactions to take place rapidly and specifically under mild temperature, relatively low concentration, and neutral conditions (*i.e.*, neither acidic nor basic). Proteins are assembled from some 20 amino acids, and, just as the 26 letters of the alphabet can be assembled in specific ways to form words of various lengths and meanings, so may tens or even hundreds of the 20 amino-acid "letters" be joined to form specific proteins. Moreover, those portions of protein molecules involved in performing similar functions in different organisms often comprise the same sequences of amino acids.

Proteins

There is the same unity among cells of all types in the manner in which living organisms preserve their individuality and transmit it to their offspring. For example, hereditary information is encoded in a specific sequence of bases that make up the DNA (deoxyribonucleic acid) molecule in the nucleus of each cell. Only four bases are used in synthesizing DNA: adenine, guanine, cytosine, and thymine. Just as the Morse Code consists of three simple signals—a dash, a dot, and a space—the precise arrangement of which suffices to convey coded messages, so the precise arrangement of the bases in DNA contains and conveys the information for the synthesis and assembly of cell components. Some primitive life-forms, however, use RNA (ribonucleic acid; a nucleic acid differing from DNA in containing the sugar ribose instead of the sugar deoxyribose and the base uracil instead of the base thymine) in place of DNA as a primary carrier of genetic information. The replication of the genetic material in these organisms must, however, pass through a DNA phase. With minor exceptions, the genetic code used by all living organisms is the same.

The chemical reactions that take place in living cells are similar as well. Green plants use the energy of sunlight

to convert water (H₂O) and carbon dioxide (CO₂) to carbohydrates (sugars and starches), other organic (carbon-containing) compounds, and molecular oxygen (O₂). The process of photosynthesis requires energy, in the form of sunlight, to split one water molecule into one-half of an oxygen molecule (O₂; the oxidizing agent) and two hydrogen atoms (H; the reducing agent), each of which dissociates to one hydrogen ion (H⁺) and one electron. Through a series of oxidation-reduction reactions, electrons (denoted e⁻) are transferred from a donating molecule (oxidation), in this case water, to an accepting molecule (reduction) by a series of chemical reactions; this "reducing power" may be coupled ultimately to the reduction of carbon dioxide to the level of carbohydrate. In effect, carbon dioxide accepts and bonds with hydrogen, forming carbohydrates (C_n[H₂O]_n).

Living organisms that require oxygen reverse this process: they consume carbohydrates and other organic materials, using oxygen synthesized by plants to form water, carbon dioxide, and energy. The process that removes hydrogen atoms (containing electrons) from the carbohydrates and passes them to the oxygen is an energy-yielding series of reactions.

In plants, all but two of the steps in the process that converts carbon dioxide to carbohydrates are the same as those steps that synthesize sugars from simpler starting materials in animals, fungi, and bacteria. Similarly, the series of reactions that take a given starting material and synthesize certain molecules that will be used in other synthetic pathways are similar, or identical, among all cell types. From a metabolic point of view, the cellular processes that take place in a lion are only marginally different from those that take place in a dandelion.

BIOLOGICAL ENERGY EXCHANGES

The energy changes associated with physicochemical processes are the province of thermodynamics, a subdiscipline of physics. The first two laws of thermodynamics state, in essence, that energy can be neither created nor destroyed and that the effect of physical and chemical changes is to increase the disorder, or randomness (*i.e.*, entropy), of the universe. Although it might be supposed that biological processes—through which organisms grow in a highly ordered and complex manner, maintain order and complexity throughout their life, and pass on the instructions for order to succeeding generations—are in contravention of these laws, this is not so. Living organisms neither consume nor create energy: they can only transform it from one form to another. From the environment they absorb energy in a form useful to them; to the environment they return an equivalent amount of energy in a biologically less useful form. The useful energy, or free energy, may be defined as energy capable of doing work under isothermal conditions (conditions in which no temperature differential exists); free energy is associated with any chemical change. Energy less useful than free energy is returned to the environment, usually as heat. Heat cannot perform work in biological systems because all parts of cells have essentially the same temperature and pressure.

The carrier of chemical energy. At any given time, a neutral molecule of water dissociates into a hydrogen ion (H⁺) and a hydroxide ion (OH⁻), and the ions are continually re-forming into the neutral molecule. Under normal conditions (neutrality), the concentration of hydrogen ions (acidic ions) is equal to that of the hydroxide ions (basic ions); each are at a concentration of 10⁻⁷ moles per litre, which is described as a pH of 7.

All cells either are bounded by membranes or contain organelles that have membranes. These membranes do not permit water or the ions derived from water to pass into or out of the cells or organelles. In green plants, sunlight is absorbed by chlorophyll and other pigments in the chloroplasts of the cells, called photosystem II. As shown previously, when a water molecule is split by light energy, one-half of an oxygen molecule and two hydrogen atoms (which dissociate to two electrons and two hydrogen ions, H⁺) are formed. When excited by sunlight, chlorophyll loses one electron to an electron carrier molecule but quickly recovers it from a hydrogen atom of the split wa-

ter molecule, which sends H⁺ into solution in the process. Two oxygen atoms come together to form a molecule of oxygen gas (O₂). The free electrons are passed to photosystem I, but, in doing so, an excess concentration of positively charged hydrogen ions (H⁺) appears on one side of the membrane in the chloroplast, whereas an excess of negatively charged hydroxide ions (OH⁻) builds up on the other side. The free energy released as H⁺ ions move through a specific "pore" in the membrane, to equalize the concentrations of ions, is sufficient to make some biological processes work, such as the uptake of certain nutrients by bacteria and the rotation of the whiplike protein-based propellers that enable such bacteria to move. Equally important, however, is that this gradient across the membrane powers the formation of adenosine triphosphate (ATP) from inorganic phosphate (HPO₄²⁻, abbreviated P_i) and adenosine diphosphate (ADP). It is ATP (Figure 1) that is the major carrier of biologically utilizable energy in all forms of living matter. The interrelationships of energy-yielding and energy-requiring metabolic reactions may be considered largely as processes that couple the formation of ATP with its breakdown.

Synthesis of ATP by green plants is similar to the synthesis of ATP that takes place in the mitochondria of animal, plant, and fungus cells, and in the plasma membranes of bacteria that use oxygen (or other inorganic electron acceptors, such as nitrate) to accept electrons from the removal of hydrogen atoms from a molecule of food (see below *The combustion of food materials: Biological energy transduction*). Through these processes most of the energy stored in food materials is released and converted into the molecules that fuel life processes. It must also be remembered, however, that many living organisms (usually

Plant synthesis of ATP

Thermodynamics

Free energy and heat

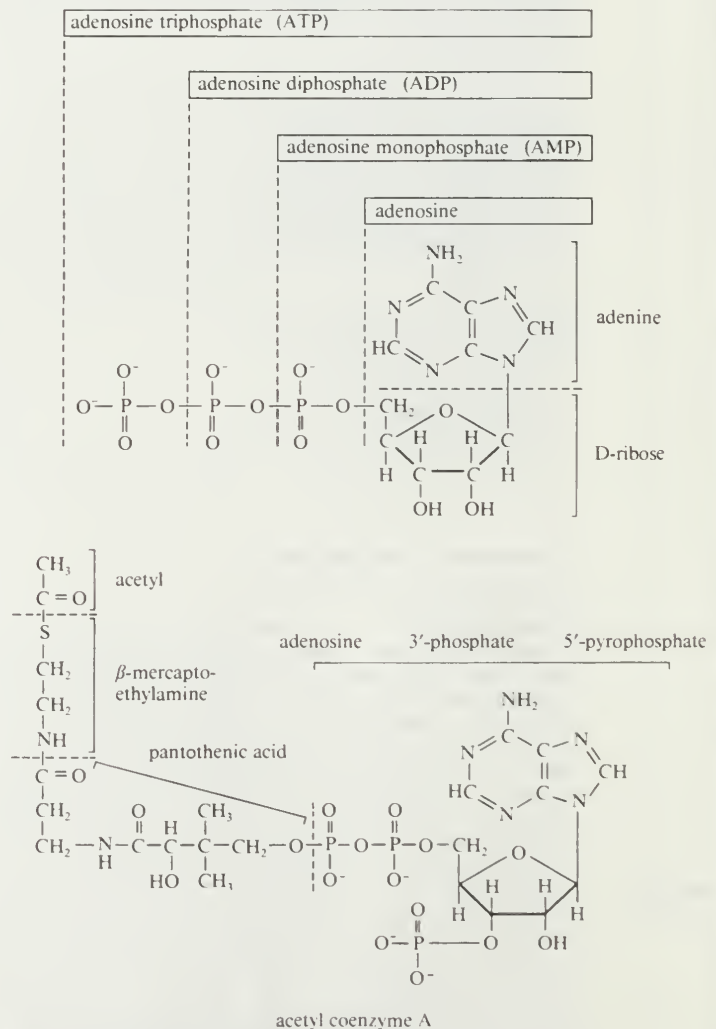


Figure 1: Biological energy carriers. (Top) ATP, ADP, and AMP. (Bottom) Acetyl coenzyme A (acetyl CoA).

bacteria and protozoa) cannot tolerate oxygen; they form ATP from inorganic phosphate and ADP by substrate-level phosphorylations (the addition of a phosphate group) that do not involve the establishment and collapse of proton gradients across membranes. Such processes are discussed in detail below (*The fragmentation of complex molecules: The catabolism of glucose*). It must also be borne in mind that the fuels of life and the cellular "furnace" in which they are "burned" are made of the same types of material: if the fires burn too brightly, not only the fuel but also the furnace is consumed. It is therefore essential to release energy at small, discrete, readily utilizable intervals. The relative complexity of the catabolic pathways (by which food materials are broken down) and the complexity of the anabolic pathways (by which cell components are synthesized) reflect this need and offer the possibility for simple feedback systems to control the rate at which materials travel along these sequences of enzymic reactions.

Catabolism. Formation of small molecules. The release of chemical energy from food materials essentially occurs in three phases. In the first phase (phase I), the large molecules that make up the bulk of food materials are broken down into small constituent units: proteins are converted to the 20 or so different amino acids of which they are composed; carbohydrates (polysaccharides such as starch in plants and glycogen in animals) are degraded to sugars such as glucose; and fats (lipids) are broken down into fatty acids and glycerol. The amounts of energy liberated in phase I are relatively small: only about 0.6 percent of the free, or useful, energy of proteins and carbohydrates, and about 0.1 percent of that of fats, is released during this phase. Because this energy is liberated largely as heat, it cannot be utilized by the cell. The purpose of the reactions of phase I, which can be grouped under the term digestion and which, in animals, occur mainly in the intestinal tract and in tissues in which reserve materials are prepared, or mobilized, for energy production, is to prepare the foodstuffs for the energy-releasing processes.

Incomplete oxidation. In the second phase of the release of energy from food (phase II), the small molecules produced in the first phase—sugars, glycerol, a number of fatty acids, and about 20 varieties of amino acids—are incompletely oxidized (in this sense, oxidation means the removal of electrons or hydrogen atoms), the end product being (apart from carbon dioxide and water) one of only three possible substances: the two-carbon compound acetate, in the form of a compound called acetyl coenzyme A (Figure 1); the four-carbon compound oxaloacetate; and the five-carbon compound α -oxoglutarate. The first, acetate in the form of acetyl coenzyme A, constitutes by far the most common product—it is the product of two-thirds of the carbon incorporated into carbohydrates and glycerol; all of the carbon in most fatty acids; and approximately half of the carbon in amino acids. The end product of several amino acids is α -oxoglutarate; that of a few others is oxaloacetate, which is formed either directly or indirectly (from fumarate). These processes, represented diagrammatically in Figure 2, show what happens in the bacterium *Escherichia coli*, but essentially similar processes occur in animals, plants, fungi, and other organisms capable of oxidizing their food materials wholly to carbon dioxide and water.

Complete oxidation. Total oxidation of the relatively few products of phase II occurs in a cyclic sequence of chemical reactions known as the tricarboxylic acid (TCA) cycle, or the Krebs cycle, after its discoverer, Sir Hans Krebs; it represents phase III of energy release from foods. Each turn of this cycle (see below *The tricarboxylic acid [TCA] cycle*) is initiated by the formation of citrate, with six carbon atoms, from oxaloacetate (with four carbons) and acetyl coenzyme A; subsequent reactions result in the reformation of oxaloacetate and the formation of two molecules of carbon dioxide. The carbon atoms that go into the formation of carbon dioxide are no longer available to the cell. The concomitant stepwise oxidations—in which hydrogen atoms or electrons are removed from intermediate compounds formed during the cycle and, via a system of carriers, are transferred ultimately to oxygen to form water—are quantitatively the most important means

of generating ATP from ADP and inorganic phosphate. These events are known as terminal respiration and oxidative phosphorylation (for details of this process, see below *Biological energy transduction*).

Some microorganisms, incapable of completely converting their carbon compounds to carbon dioxide, release energy by fermentation reactions, in which the intermediate compounds of catabolic routes either directly or indirectly accept or donate hydrogen atoms. Such secondary changes in intermediate compounds result in considerably less energy being made available to the cell than occurs with the pathways that are linked to oxidative phosphorylation; however, fermentation reactions yield a large variety of commercially important products. Thus, for example, if the oxidation (removal of electrons or hydrogen atoms) of some catabolic intermediate is coupled to the reduction of pyruvate or of acetaldehyde derived from pyruvate, the products formed are lactic acid and ethyl alcohol, respectively.

Anabolism. Catabolic pathways effect the transformation of food materials into the interconvertible intermediates of the pathways shown in Figure 2. Anabolic pathways, on the other hand, are sequences of enzyme-catalyzed reactions in which the component building blocks of large molecules, or macromolecules (e.g., proteins, carbohydrates, and fats), are constructed from the same intermediates. Thus, catabolic routes have clearly defined beginnings but no unambiguously identifiable end products; anabolic routes, on the other hand, lead to clearly distinguishable end products from diffuse beginnings. The two types of pathway are linked through reactions of phosphate transfer, involving ADP, AMP, and ATP as described above, and also through electron transfers, which enable reducing equivalents (*i.e.*, hydrogen atoms or electrons), which have been released during catabolic reactions, to be utilized for biosynthesis. But, although catabolic and anabolic pathways are closely linked, and although the overall effect of one type of route is obviously the opposite of the other, they have few steps in common. The anabolic pathway for the synthesis of a particular molecule generally starts from intermediate compounds quite different from those produced as a result of catabolism of that molecule; for example, microorganisms catabolize aromatic (*i.e.*, containing a ring, or cyclic, structure) amino acids to acetyl coenzyme A and an intermediate compound of the TCA cycle. The biosynthesis of these amino acids, however, starts with a compound derived from pyruvate and an

Anabolism versus catabolism

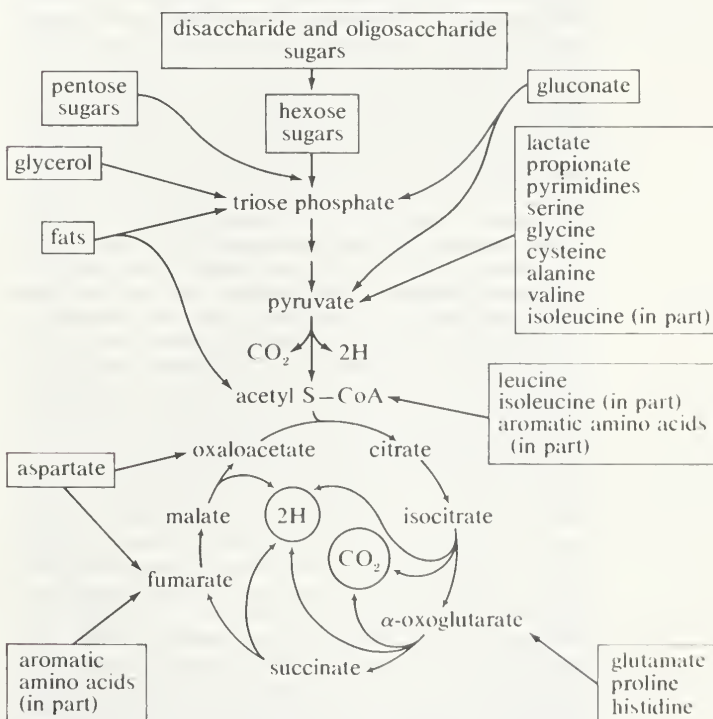


Figure 2: Pathways for the catabolism of nutrients by *Escherichia coli*.

Phase I

Phase II

Phase III

intermediate compound of the metabolism of pentose (a general name for sugars with five carbon atoms). Similarly, histidine is synthesized from a pentose sugar but is catabolized to α -oxoglutarate.

Even in cases in which a product of catabolism is used in an anabolic pathway, differences emerge; thus, fatty acids, which are catabolized to acetyl coenzyme A, are synthesized not from acetyl coenzyme A directly but from a derivative of it, malonyl coenzyme A (see below *The biosynthesis of cell components: Lipid components*). Furthermore, even enzymes that catalyze apparently identical steps in catabolic and anabolic routes may exhibit different properties. In general, therefore, the way down (catabolism) is different from the way up (anabolism). These differences are important because they allow for the regulation of catabolic and anabolic processes in the cell.

In eukaryotic cells (*i.e.*, those with a well-defined nucleus, characteristic of organisms higher than bacteria) the enzymes of catabolic and anabolic pathways are often located in different cellular compartments. This also contributes to the manner of their cellular control; for example, the formation of acetyl coenzyme A from fatty acids, referred to above, occurs in animal cells in small sausage-shaped components, or organelles, called mitochondria, which also contain the enzymes for terminal respiration and for oxidative phosphorylation. The biosynthesis of fatty acids from acetyl coenzyme A, on the other hand, occurs in the cytoplasm.

INTEGRATION OF CATABOLISM AND ANABOLISM

Fine control. Possibly the most important means for controlling the flux of metabolites through catabolic and anabolic pathways, and for integrating the numerous different pathways in the cell, is through the regulation of either the activity or the synthesis of key (pacemaker) enzymes. It was recognized in the 1950s, largely from work with microorganisms, that pacemaker enzymes can interact with small molecules at more than one site on the surface of the enzyme molecule. The reaction between an enzyme and its substrate—defined as the compound with which the enzyme acts to form a product—occurs at a specific site on the enzyme known as the catalytic, or active, site; the proper fit between the substrate and the active site is an essential prerequisite for the occurrence of a reaction catalyzed by an enzyme (see *BIOCHEMICAL COMPONENTS OF ORGANISMS*). Interactions at other, so-called regulatory sites on the enzyme, however, do not result in a chemical reaction but cause changes in the shape of the protein; the changes profoundly affect the catalytic properties of the enzyme, either inhibiting or stimulating the rate of the reaction. Modulation of the activity of pacemaker enzymes may be effected by metabolites of the pathway in which the enzyme acts or by those of another pathway; the process may be described as a “fine control” of metabolism. Very small changes in the chemical environment thus produce important and immediate effects on the rates at which individual metabolic processes occur.

Most catabolic pathways are regulated by the relative proportions of ATP, ADP, and AMP in the cell. It is reasonable to suppose that a pathway that serves to make ATP available for energy-requiring reactions would be less active if sufficient ATP were already present, than if ADP or AMP were to accumulate. The relative amounts of the adenine nucleotides (*i.e.*, ATP, ADP, and AMP) thus modulate the overall rate of catabolic pathways. They do so by reacting with specific regulatory sites on pacemaker enzymes necessary for the catabolic pathways, which do not participate in the anabolic routes that effect the opposite reactions. Similarly, it is reasonable to suppose that many anabolic processes, which require energy, are inhibited by ADP or AMP; elevated levels of these nucleotides may be regarded therefore as cellular distress signals indicating a lack of energy.

Since one way in which anabolic pathways differ from catabolic routes is that the former result in identifiable end products, it is not unexpected that the pacemaker enzymes of many anabolic pathways—particularly those effecting the biosynthesis of amino acids and nucleotides—are regulated by the end products of these pathways

or, in cases in which branching of pathways occurs, by end products of each branch. Such pacemaker enzymes usually act at the first step unique to a particular anabolic route. If branching occurs, the first step of each branch is controlled. By this so-called negative feedback system, the cellular concentrations of products determine the rates of their formation, thus ensuring that the cell synthesizes only as much of the products as it needs.

Coarse control. A second and less immediately responsive, or “coarse,” control is exerted over the synthesis of pacemaker enzymes. The rate of protein synthesis reflects the activity of appropriate genes, which contain the information that directs all cellular processes. Coarse control is therefore exerted on genetic material rather than on enzymes. Preferential synthesis of a pacemaker enzyme is particularly required to accommodate a cell to major changes in its chemical milieu. Such changes occur in multicellular organisms only to a minor extent, so that this type of control mechanism is less important in animals than in microorganisms. In the latter, however, it may determine the ease with which a cell previously growing in one nutrient medium can grow after transfer to another. In cases in which several types of organism compete in the same medium for available carbon sources, the operation of coarse controls may well be decisive in ensuring survival.

Alterations in the differential rates of synthesis of pacemaker enzymes in microorganisms responding to changes in the composition of their growth medium also manifest the properties of negative feedback systems. Depending on the nature of the metabolic pathway of which a pacemaker enzyme is a constituent, the manner in which the alterations are elicited may be distinguished. Thus, an increase in the rates at which enzymes of catabolic routes are synthesized results from the addition of inducers—usually compounds that exhibit some structural similarity to the substrates on which the enzymes act. A classic example of an inducible enzyme of this type is β -galactosidase. *Escherichia coli* growing in nutrient medium containing glucose do not utilize the milk sugar, lactose (glucose-4- β -D-galactoside); however, if the bacteria are placed in a growth medium containing lactose as the sole source of carbon, they synthesize β -galactosidase and can therefore utilize lactose. The reaction catalyzed by the enzyme is the hydrolysis (*i.e.*, breakdown involving water) of lactose to its two constituent sugars, glucose and galactose; the preferential synthesis of the enzyme thus allows the bacteria to use the lactose for growth and energy. Another characteristic of the process of enzyme induction is that it continues only as long as the inducer (in this case, lactose) is present; if cells synthesizing β -galactosidase are transferred to a medium containing no lactose, synthesis of β -galactosidase ceases, and the amount of the enzyme in the cells is diluted as they divide, until the original low level of the enzyme is reestablished.

In contrast, the differential rates of synthesis of pacemaker enzymes of anabolic routes are usually not increased by the presence of inducers. Instead, the absence of small molecules that act to repress enzyme synthesis accelerates enzyme formation. Similar to the fine control processes described above is the regulation by coarse control of many pacemaker enzymes of amino-acid biosynthesis. Like the end product inhibitors, the repressors in these cases also appear to be the amino-acid end products themselves.

It is useful to regard the acceleration of the enzyme-forming machinery as the consequence, metaphorically, of either placing a foot on the accelerator or removing it from the brake. Analysis of the mechanisms by which gene activity is controlled suggest, however, that the distinction between inducible and repressible enzymes may be more apparent than real (see below *Regulation of metabolism*).

THE STUDY OF METABOLIC PATHWAYS

There are two main reasons for studying a metabolic pathway: (1) to describe, in quantitative terms, the chemical changes catalyzed by the component enzymes of the route; and (2) to describe the various intracellular controls that govern the rate at which the pathway functions.

Studies with whole organisms or organs can provide

Eukaryotic cells

Gene involvement

Purpose

information that one substance is converted to another and that this process is localized in a certain tissue; for example, experiments can show that urea, the chief nitrogen-containing end product of protein metabolism in mammals, is formed exclusively in the liver. They cannot reveal, however, the details of the enzymatic steps involved. Clues to the identity of the products involved, and to the possible chemical changes effected by component enzymes, can be provided in any of four ways involving studies with either whole organisms or tissues.

First, under stress or the imbalances associated with diseases, certain metabolites may accumulate to a greater extent than normal. Thus, during the stress of violent exercise, lactic acid appears in the blood, while glycogen, the form in which carbohydrate is stored in muscle, disappears. Such observations do not, however, prove that lactic acid is a normal intermediate of glycogen catabolism; rather, they show only that compounds capable of yielding lactic acid are likely to be normal intermediates. Indeed, in the example, lactic acid is formed in response to abnormal circumstances and is not directly formed in the pathways of carbohydrate catabolism. On the other hand, the abnormal accumulation of pyruvic acid in the blood of vitamin B₁-deficient pigeons was a valuable clue to the role of this vitamin in the oxidation of pyruvate.

Second, the administration of metabolic poisons may lead to the accumulation of specific metabolites. If fluoroacetic acid or fluorocitric acid is ingested by animals, for example, citric acid accumulates in the liver. This correctly suggests that fluorocitric acid administered as such, or formed from fluoroacetic acid via the tricarboxylic acid (TCA) cycle, inhibits an enzyme of citrate oxidation.

Third, the fate of any nutrient—indeed, often the fate of a particular chemical group or atom in a nutrient—can be followed with relative ease by administering the nutrient labeled with an isotope. Isotopes are forms of an element that are chemically indistinguishable from each other but differ in physical properties.

The use of a nonradioactive isotope of nitrogen in the 1930s first revealed the dynamic state of body constituents. It had previously been believed that the proteins of tissues are stable once formed, disappearing only with the death of the cell. By feeding amino acids labeled with isotopic nitrogen to rats, it was discovered that the isotope was incorporated into many of the amino acids found in proteins of the liver and the gut, even though the total protein content of these tissues did not change. This suggested that the proteins of these tissues exist in a dynamic steady state, in which relatively high rates of synthesis are counterbalanced by equal rates of degradation. Thus, although the average liver cell has a life-span of several months, half of its proteins are synthesized and degraded every five to six days. On the other hand, the proteins of the muscle or the brain, tissues that (unlike the gut or liver) need not adjust to changes in the chemical composition of their milieu, do not turn over as rapidly. The high rates of turnover observed in liver and gut tissues indicate that the coarse controls, exerted through the onset and cessation of synthesis of pacemaker enzymes, do occur in animal cells.

Finally, genetically altered organisms (mutants) fail to synthesize certain enzymes in an active form. Such defects, if not lethal, result in the accumulation and excretion of the substrate of the defective enzyme; in normal organisms, the substrate would not accumulate, because it would be acted upon by the enzyme. The significance of this observation was first realized in the early 20th century when the phrase "inborn errors of metabolism" was used to describe hereditary conditions in which a variety of amino acids and other metabolites are excreted in the urine (see below *Metabolic diseases*). In microorganisms, in which it is relatively easy to cause genetic mutations and to select specific mutants, this technique has been very useful. In addition to their utility in the unraveling of metabolic pathways, the use of mutants in the early 1940s led to the postulation of the one gene-one enzyme hypothesis by the Nobel Prize winners George W. Beadle and Edward L. Tatum; their discoveries opened the field of biochemical genetics and first revealed the nature of the fine controls of metabolism.

Because detailed information about the mechanisms of component enzymatic steps in any metabolic pathway cannot be obtained from studies with whole organisms or tissues, various techniques have been developed for studying these processes—*e.g.*, sliced tissues, and homogenates and cell-free extracts, which are produced by physical disruption of the cells and the removal of cell walls and other debris. The sliced-tissue technique was successfully used by the Nobel Prize winner Sir Hans Krebs in his pioneer studies in the early 1930s on the mechanism of urea formation in the liver. Measurements were made of the stimulating effects of small quantities of amino acids on both the rate of oxygen uptake and the amount of oxygen taken up; the amino acids were added to liver slices bathed in a nutrient medium. Such measurements revealed the cyclic nature of the process; specific amino acids acted as catalysts, stimulating respiration to an extent greater than expected from the quantities added. This was because the added material had been re-formed in the course of the cycle (see below *The catabolism of proteins: Disposal of nitrogen*).

Homogenates of tissue are useful in studying metabolic processes because permeability barriers that may prevent ready access of external materials to cell components are destroyed. The tissue is usually minced, blended, or otherwise disrupted in a medium that is suitably buffered to maintain the normal acid-base balance of the tissue, and contains the ions required for many life processes, chiefly sodium, potassium, and magnesium. The tissue is either used directly—as was done by Krebs in elucidating, in 1937, the TCA cycle from studies of the respiration of minced pigeon breast muscle—or fractionated (*i.e.*, broken down) further. If the latter procedure is followed, homogenization is often carried out in a medium containing a high concentration of the sugar sucrose, which provides a milieu favourable for maintaining the integrity of cellular components. The components are recovered by careful spinning in a centrifuge, at a series of increasing speeds. It is thus possible to obtain fractions containing predominantly one type of organelle: nuclei (and some unbroken cells); mitochondria, lysosomes, and microbodies; microsomes (*i.e.*, ribosomes and endoplasmic reticulum fragments); and—after prolonged centrifugation at forces in excess of 100,000 times gravity—a clear liquid that represents the soluble fraction of the cytoplasm. The fractions thus obtained can be further purified and tested for their capacity to carry out a given metabolic step or steps. This procedure was used to show that isolated mitochondria catalyze the oxidation reactions of the TCA cycle and that these organelles also contain the enzymes of fatty acid oxidation. Similarly, isolated ribosomes are used to study the pathway and mechanism of protein synthesis.

The final step in elucidating a reaction in a metabolic pathway includes isolation of the enzyme involved. The rate of the reaction and the factors that control the activity of the enzyme are then measured.

It should be emphasized that biochemists realize that studies on isolated and highly purified systems, such as those briefly described above, can do no more than approximate biological reality. The identification of the fine and coarse controls of a metabolic pathway, and (when appropriate) other influences on that pathway, must ultimately involve the study of the pathway in the whole cell or organism. Although some techniques have proved adequate for relating findings in the test tube to the situation in living organisms, study of the more complex metabolic processes, such as those involved in differentiation and development, may require the elaboration of new experimental approaches.

The fragmentation of complex molecules

Food materials must undergo oxidation in order to yield biologically useful energy. Oxidation does not necessarily involve oxygen, although it must involve the transfer of electrons from a donor molecule to a suitable acceptor molecule; the donor is thus oxidized and the recipient reduced. Many microorganisms either must live in the absence of oxygen (*i.e.*, are obligate anaerobes) or can live in

The use of homogenates

Identification of metabolites

Rate of protein turnover

its presence or its absence (*i.e.*, are facultative anaerobes).

Fermentations

If no oxygen is available, the catabolism of food materials is effected via fermentations, in which the final acceptor of the electrons removed from the nutrient is some organic molecule, usually generated during the fermentation process. There is no net oxidation of the food molecule in this type of catabolism; that is, the overall oxidation state of the fermentation products is the same as that of the starting material.

Organisms that can use oxygen as a final electron acceptor also use many of the steps in the fermentation pathways in which food molecules are broken down to smaller fragments; these fragments, instead of serving as electron acceptors, are fed into the TCA cycle, the pathway of terminal respiration.

In this cycle all of the hydrogen atoms (H) or electrons (e^-) are removed from the fragments and are channeled through a series of electron carriers, ultimately to react with oxygen (O; see below *Energy conservation*). All carbon atoms are eliminated as carbon dioxide (CO_2) in this process. The sequence of reactions involved in the catabolism of food materials may thus be conveniently considered in terms of an initial fragmentation (fermentation), followed by a combustion (respiration) process.

THE CATABOLISM OF GLUCOSE

Glycolysis. *The transformation of glucose.* Quantitatively, the most important source of energy for cellular processes is the six-carbon sugar glucose ($\text{C}_6\text{H}_{12}\text{O}_6$). Two structures of glucose are shown in Figure 3, in which the carbon atoms are numbered. (See *BIOCHEMICAL COMPONENTS OF ORGANISMS* for a discussion of the chemical nature of glucose and other carbohydrates.) Glucose is made available to animals through the hydrolysis of polysaccharides, such as glycogen and starch, the process being catalyzed by digestive enzymes. In animals, the sugar thus set free passes from the gut into the bloodstream and from there into the cells of the liver and other tissues. In microorganisms, of course, no such specialized tissues are involved.

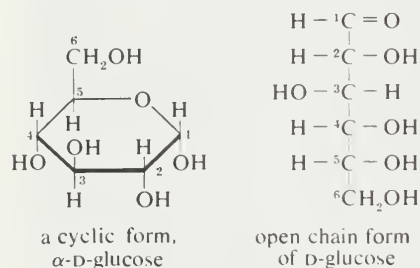
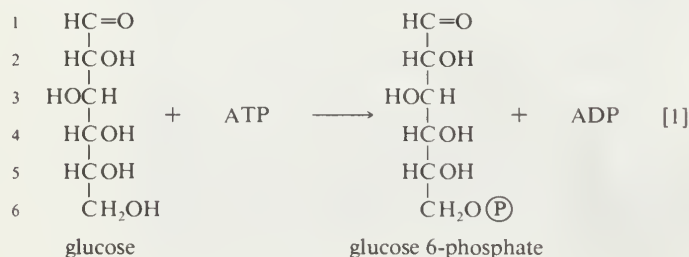


Figure 3: The structure of α -D-glucose.

The fermentative phase of glucose catabolism (glycolysis) involves several enzymes; the action of each is summarized below. In living cells many of the compounds that take part in metabolism exist as negatively charged moieties, or anions, and are named as such in most of this article; *e.g.*, pyruvate, oxaloacetate.

In order to obtain a net yield of ATP from the catabolism of glucose, it is first necessary to invest ATP. During step [1] the alcohol group at position 6 of the glucose molecule readily reacts with the terminal phosphate group of ATP, forming glucose 6-phosphate and ADP. For convenience, the phosphoryl group (PO_3^{2-}) is represented by P . Because the decrease in free energy is so large, this reaction is virtually irreversible under physiological conditions.

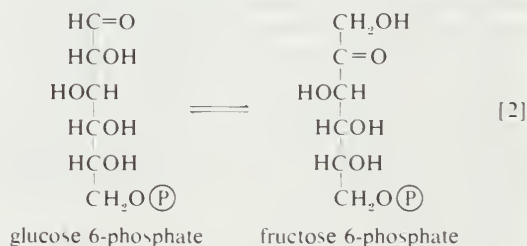
The formation of glucose 6-phosphate



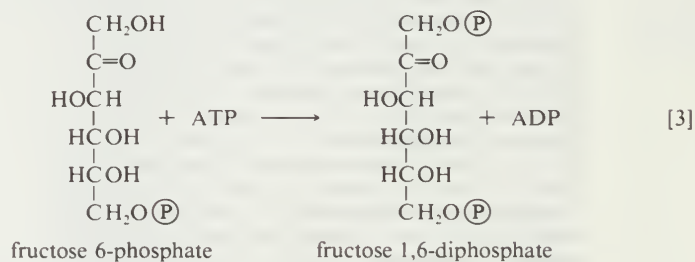
In animals, this phosphorylation of glucose, which yields glucose 6-phosphate, is catalyzed by two different enzymes. In most cells a hexokinase with a high affinity for glucose—*i.e.*, only small amounts of glucose are necessary for enzymatic activity—effects the reaction. In addition, the liver contains a glucokinase, which requires a much greater concentration of glucose before it reacts. Glucokinase functions only in emergencies, when the concentration of glucose in the blood rises to abnormally high levels.

Certain facultative anaerobic bacteria also contain hexokinases but apparently do not use them to phosphorylate glucose. In such cells, external glucose can be utilized only if it is first phosphorylated to glucose 6-phosphate via a system linked to the cell membrane that involves a compound called phosphoenolpyruvate (formed in step [9] of glycolysis), which serves as an obligatory donor of the phosphate group; *i.e.*, ATP cannot serve as the phosphate donor in the reaction.

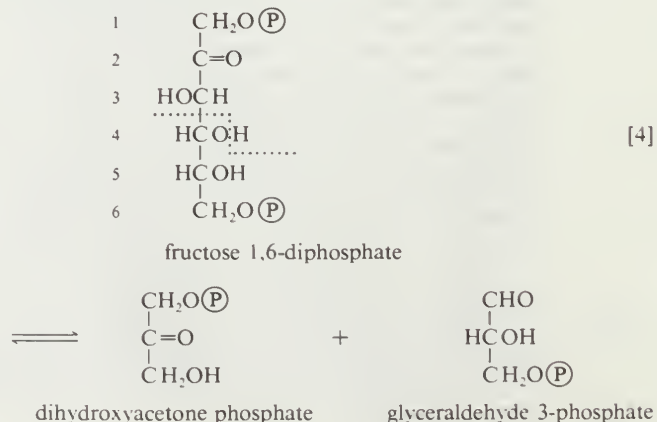
The reaction in which glucose 6-phosphate is changed to fructose 6-phosphate is catalyzed by phosphoglucose isomerase [2]. In the reaction, a secondary alcohol group ($-\text{CHOH}$) at the second carbon atom is oxidized to a keto-group (*i.e.*, $-\text{C}=\text{O}$), and the aldehyde group ($-\text{CHO}$) at the first carbon atom is reduced to a primary alcohol group ($-\text{CH}_2\text{OH}$). Reaction [2] is readily reversible, as is indicated by the double arrows.



The formation of the alcohol group at the first carbon atom permits the repetition of the reaction effected in step [1]; that is, a second molecule of ATP is invested. The product is fructose 1,6-diphosphate [3]. Again, as in the hexokinase reaction, the decrease in free energy of the reaction, which is catalyzed by phosphofructokinase, is sufficiently large to make this reaction virtually irreversible under physiological conditions: ADP is also a product.

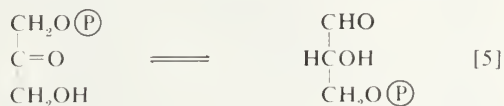


The first three steps of glycolysis have thus transformed an asymmetrical sugar molecule, glucose, into a symmetri-



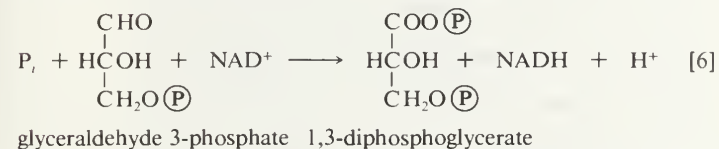
cal form, fructose 1,6-diphosphate, containing a phosphoryl group at each end; the molecule next is split into two smaller fragments that are interconvertible. This elegant simplification is achieved via steps [4] and [5], which are described below.

The aldolase reaction. In [4], an enzyme catalyzes the breaking apart of the six-carbon sugar fructose 1,6-diphosphate into two three-carbon fragments. The molecule is split between carbons 3 and 4. Reversal of this cleavage—*i.e.*, the formation of a six-carbon compound from two three-carbon compounds—is possible. Because the reverse reaction is an aldol condensation—*i.e.*, an aldehyde (glyceraldehyde 3-phosphate) combines with a ketone (dihydroxyacetone phosphate)—the enzyme is commonly called aldolase. The two three-carbon fragments produced in step [4], dihydroxyacetone phosphate and glyceraldehyde 3-phosphate, are also called triose phosphates. They are readily converted to each other by a process [5] analogous to that in step [2]. The enzyme that catalyzes the interconversion [5] is triose phosphate isomerase, a different enzyme than that catalyzing step [2].



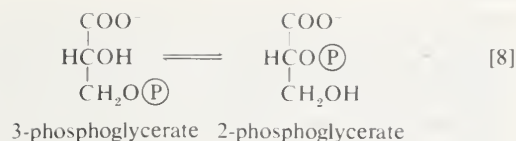
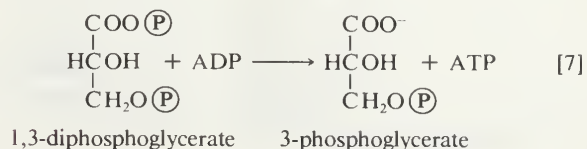
dihydroxyacetone phosphate glyceraldehyde 3-phosphate

The formation of ATP. The second stage of glucose catabolism comprises reactions [6] through [10], in which a net gain of ATP is achieved through the oxidation of one of the triose phosphate compounds formed in step [5]. One molecule of glucose forms two molecules of the triose phosphate; both three-carbon fragments follow the same pathway, and steps [6] through [10] must occur twice to complete the glucose breakdown.



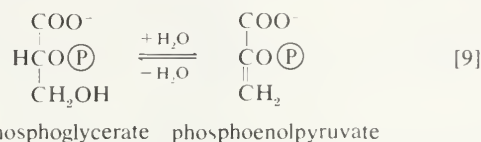
Step [6], in which glyceraldehyde 3-phosphate is oxidized, is one of the most important reactions in glycolysis. It is during this step that the energy liberated during oxidation of the aldehyde group ($-\text{CHO}$) is conserved in the form of a high-energy phosphate compound; namely, as 1,3-diphosphoglycerate, an anhydride of a carboxylic acid and phosphoric acid. The hydrogen atoms or electrons removed from the aldehyde group during its oxidation are accepted by a coenzyme (so called because it functions in conjunction with an enzyme) involved in hydrogen or electron transfer; the coenzyme, nicotinamide adenine dinucleotide (NAD^+), is reduced to form $\text{NADH} + \text{H}^+$ in the process. The NAD^+ thus reduced is bound to the enzyme glyceraldehyde 3-phosphate dehydrogenase, catalyzing the overall reaction, step [6].

The 1,3-diphosphoglycerate produced in step [6] reacts with ADP in a reaction catalyzed by phosphoglycerate kinase, with the result that one of the two phosphoryl groups is transferred to ADP to form ATP and 3-phosphoglycerate. This reaction [7] is highly exergonic (*i.e.*, it proceeds with a loss of free energy); as a result, the oxidation of glyceraldehyde 3-phosphate, step [6], is irreversible. In summary, the energy liberated during oxidation of an aldehyde group ($-\text{CHO}$ in glyceraldehyde 3-phosphate) to a carboxylic acid group ($-\text{COO}^-$ in 3-phosphoglycerate) is conserved as the phosphate bond energy in ATP during steps [6] and [7]. This step occurs twice for each molecule of glucose; thus the initial investment of ATP in steps [1] and [3] is recovered.

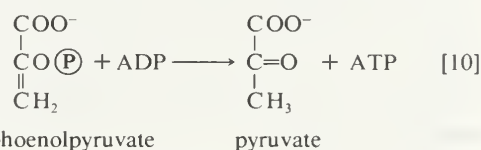


The 3-phosphoglycerate in step [7] now forms 2-phosphoglycerate, in a reaction catalyzed by phosphoglyceromutase [8]. During step [9] the enzyme enolase reacts with 2-phosphoglycerate to form phosphoenolpyruvate (PEP), water being lost from 2-phosphoglycerate in the process. Phosphoenolpyruvate acts as the second source of ATP in glycolysis. The transfer of the phosphate group from PEP to ADP, catalyzed by pyruvate kinase [10], is also highly exergonic and is thus virtually irreversible under physiological conditions.

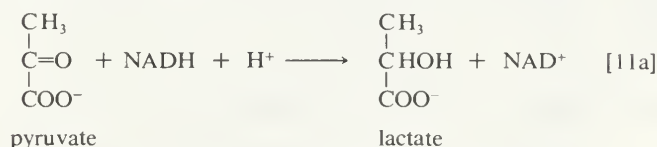
A second source of ATP



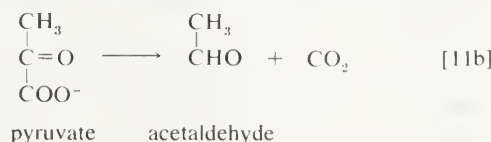
Reaction [10] occurs twice for each molecule of glucose entering the glycolytic sequence; thus the net yield is two molecules of ATP for each six-carbon sugar. No further molecules of glucose can enter the glycolytic pathway, however, until the $\text{NADH} + \text{H}^+$ produced in step [6] is reoxidized to NAD^+ . In anaerobic systems this means that electrons must be transferred from ($\text{NADH} + \text{H}^+$) to some



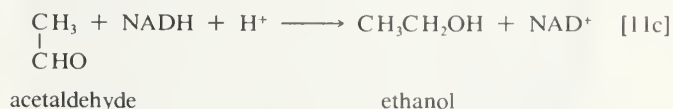
organic acceptor molecule, which thus is reduced in the process. Such an acceptor molecule could be the pyruvate formed in reaction [10]. In certain bacteria (*e.g.*, so-called lactic acid bacteria) or in muscle cells functioning vigorously in the absence of adequate supplies of oxygen,



pyruvate is reduced to lactate via a reaction catalyzed by lactate dehydrogenase (reaction [11a]); *i.e.*, NADH gives up its hydrogen atoms or electrons to pyruvate, and lactate and NAD^+ are formed. Alternatively, in organisms such as brewers' yeast, pyruvate is first decarboxylated to form



acetaldehyde and carbon dioxide in a reaction catalyzed by pyruvate decarboxylase [11b]; acetaldehyde then is reduced (by $\text{NADH} + \text{H}^+$) in a reaction catalyzed by alcohol dehydrogenase [11c], yielding ethanol and oxidized coenzyme (NAD^+).



Many variations of reaction [11] occur in nature. In the heterolactic (mixed lactic acid) fermentations carried out by some microorganisms, a mixture of reactions [11a, b, and c] regenerates NAD^+ and results in the production, for each molecule of glucose fermented, of a molecule each of lactate, ethanol, and carbon dioxide. In other types

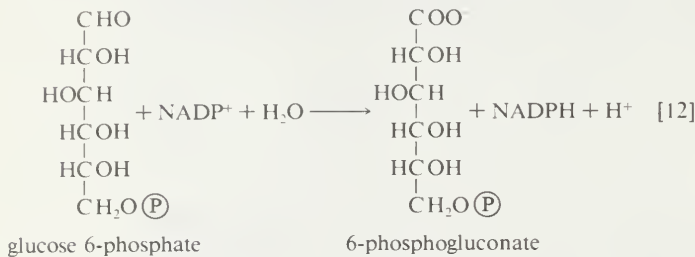
Aldol condensation

Conversion of energy of oxidation into ATP

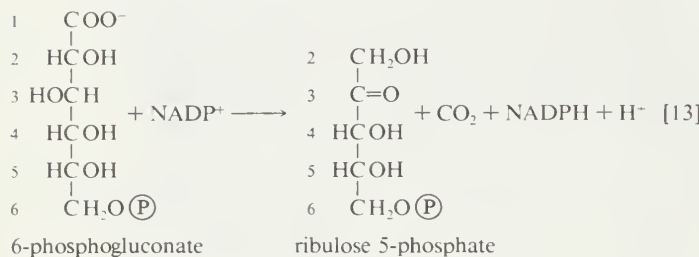
of fermentation, the end products may be derivatives of acids such as propionic, butyric, acetic, and succinic; decarboxylated materials derived from them (e.g., acetone); or compounds such as glycerol.

The phosphogluconate pathway. Many cells possess, in addition to all or part of the glycolytic pathway that comprises reactions [1] through [11], other pathways of glucose catabolism that involve, as the first unique step, the oxidation of glucose 6-phosphate [12] instead of the formation of fructose 6-phosphate [2]. This is the phosphogluconate pathway, or pentose phosphate cycle. During reaction [12], hydrogen atoms or electrons are removed from the carbon atom at position 1 of glucose 6-phosphate in a reaction catalyzed by glucose 6-phosphate dehydrogenase. The product of the reaction is 6-phosphogluconate.

The oxidation of glucose 6-phosphate



The reducing equivalents (hydrogen atoms or electrons) are accepted by nicotine adenine dinucleotide phosphate (NADP⁺), a coenzyme similar to but not identical with NAD⁺. A second molecule of NADP⁺ is reduced as 6-phosphogluconate is further oxidized; the reaction is catalyzed by 6-phosphogluconate dehydrogenase [13]. The products of the reaction also include ribulose 5-phosphate and carbon dioxide. (The numbers at the carbon atoms in step [13] indicate that carbon 1 of 6-phosphogluconate forms carbon dioxide.)



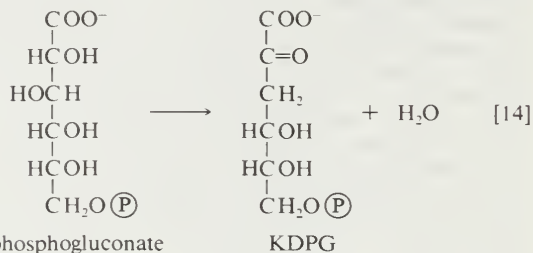
Ribulose 5-phosphate can undergo a series of reactions in which two-carbon and three-carbon fragments are interchanged between a number of sugar phosphates; this sequence of events can lead to the formation of two molecules of fructose 6-phosphate and one of glyceraldehyde 3-phosphate from three molecules of ribulose 5-phosphate (i.e., the conversion of three molecules with five carbons to two with six and one with three). Although the cycle, which is outlined in Figure 4, is the main pathway in microorganisms for fragmentation of pentose sugars, it is not of major importance as a route for the oxidation of glucose. Its primary purpose in most cells is to generate reducing power in the cytoplasm, in the form of reduced NADP⁺. This function is especially prominent in tissues—such as the liver, mammary gland, fat tissue, and the cortex (outer region) of the adrenal gland—that actively carry out the biosynthesis of fatty acids and other fatty substances (e.g., steroids). A second function of reactions [12] and [13] is to generate from glucose 6-phosphate the pentoses that are used in the synthesis of nucleic acids (see below *The biosynthesis of cell components*).

In photosynthetic organisms, some of the reactions of the phosphogluconate pathway are part of the major route for the formation of sugars from carbon dioxide; in this case, the reactions occur in a direction opposite to that in which they occur in nonphotosynthetic tissues (see PHOTOSYNTHESIS).

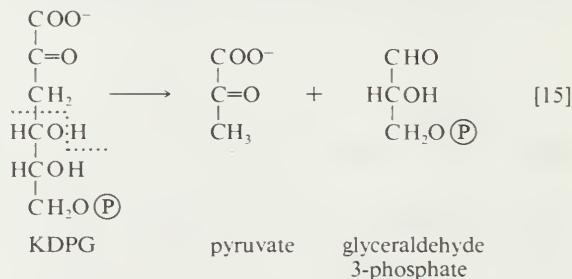
Alternative route

A different route for the catabolism of glucose also involves 6-phosphogluconate; it is of considerable importance in microorganisms lacking some of the enzymes

necessary for glycolysis. In this route, 6-phosphogluconate (derived from glucose via steps [1] and [12]) is not oxidized to ribulose 5-phosphate via reaction [13] but, in an enzyme-catalyzed reaction [14], loses water, forming the compound 2-keto-3-deoxy-6-phosphogluconate (KDPG).



This is then split into pyruvate and glyceraldehyde-3-phosphate [15], both of which are intermediates of the glycolytic pathway.



THE CATABOLISM OF SUGARS OTHER THAN GLUCOSE

The main storage carbohydrate of animal cells is glycogen, in which chains of glucose molecules—linked end-to-end, the C1 position of one glucose being linked to the C4 position of the adjacent one—are joined to each other by occasional linkages between a carbon at position 1 on one glucose and a carbon at position 6 on another. Two enzymes cooperate in releasing glucose molecules from glycogen. Glycogen phosphorylase catalyzes the splitting of the 1,4-bonds by adding the elements of phosphoric acid

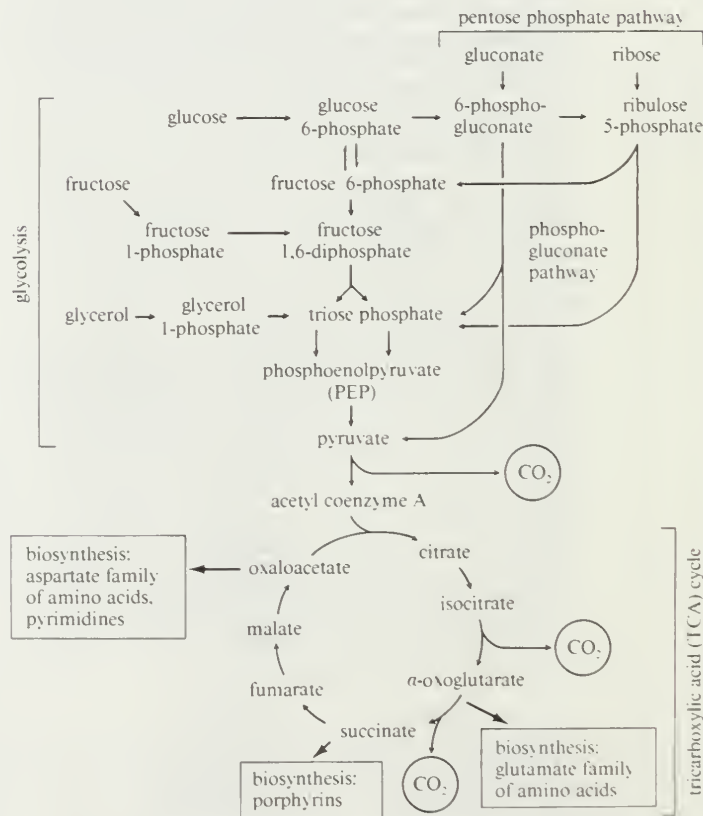
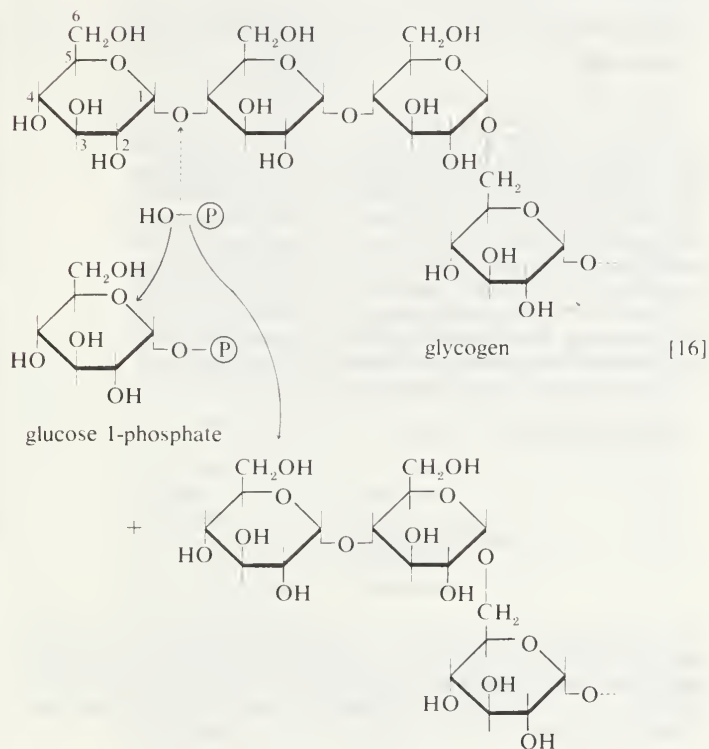


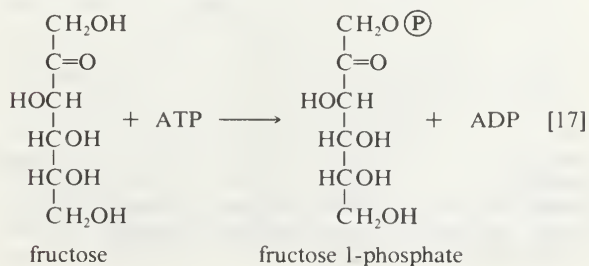
Figure 4 Pathways for the utilization of carbohydrates.

at the point shown by the broken arrow in [16], rather than water, as in the digestive hydrolysis of polysaccharides such as glycogen and starch. The products of [16] are glucose 1-phosphate and chains of sugar molecules shortened by one unit; the chains are degraded further by repetition of step [16]. When a bridge linking two chains, at C1 and C6 carbon atoms of adjacent glucose units, is reached, it is hydrolyzed in a reaction involving the enzyme α (1 \rightarrow 6) glucosidase. After the two chains are separated, reaction [16] can occur again. The glucose 1-phosphate thus formed from glycogen or, in plants, from starch, is converted to glucose 6-phosphate by phosphoglucomutase [78], which catalyzes a reaction very similar to that effected in step [8] of glycolysis; glucose 6-phosphate can then undergo further catabolism via glycolysis [2-10] or via either of the routes involving formation of 6-phosphogluconate [12].

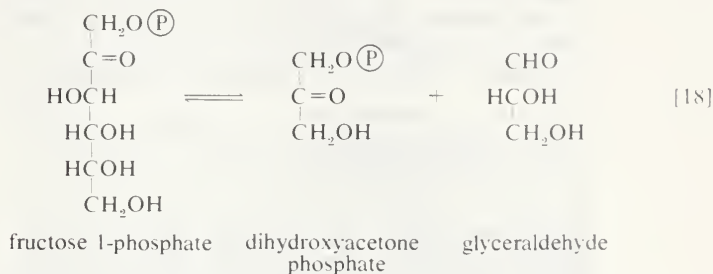


Fragmentation of lactose and sucrose

Other sugars encountered in the diet are likewise transformed to products that are intermediates of central metabolic pathways. Lactose, or milk sugar, is composed of one molecule of galactose linked to one molecule of glucose. Sucrose, the common sugar of cane or beet, is made up of glucose linked to fructose. Both sucrose and lactose are hydrolyzed to glucose and fructose or galactose, respectively. Glucose is utilized as already described, but special reactions must occur before the other sugars can enter the catabolic routes. Galactose, for example, is phosphorylated in a manner analogous to step [1] of glycolysis. The reaction, catalyzed by a galactokinase, results in the formation of galactose 1-phosphate; this product is transformed to glucose 1-phosphate by a sequence of reactions requiring as a coenzyme uridine triphosphate (UTP). Fructose may also be phosphorylated in animal cells through the action of hexokinase [1], in which case fructose 6-phosphate is the product, or in liver tissue via a fructokinase that gives rise to fructose 1-phosphate [17]. Adenosine triphosphate supplies the phosphate group in both cases.



Fructose 1-phosphate is also formed when facultative anaerobic microorganisms use fructose as a carbon source for growth; in this case, however, the source of the phosphate is phosphoenolpyruvate rather than ATP. Fructose 1-phosphate can be catabolized by one of two routes. In the liver, it is split by an aldolase enzyme [18] abundant in that tissue (but lacking in muscle); the products are dihydroxyacetone phosphate and glyceraldehyde. It will be recalled that dihydroxyacetone phosphate is an intermediate compound of glycolysis. Although glyceraldehyde is not an intermediate of glycolysis, it can be converted to one (glyceraldehyde 3-phosphate) in a reaction involving the conversion of ATP to ADP.



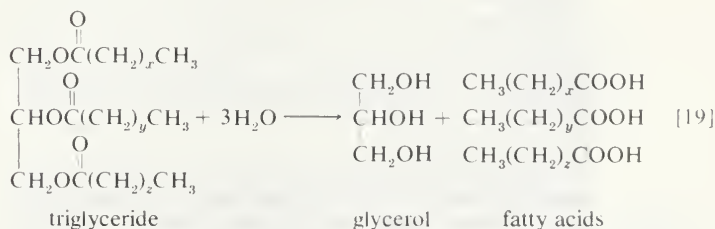
In many organisms other than mammals, fructose 1-phosphate does not have to undergo reaction [18] in order to enter central metabolic routes. Instead, a fructose 1-phosphate kinase, distinct from the phosphofructokinase that catalyzes step [3] of glycolysis, effects the direct conversion of fructose 1-phosphate and ATP to fructose 1,6-diphosphate and ADP.

THE CATABOLISM OF LIPIDS (FATS)

Although carbohydrates are the major fuel for most organisms, fatty acids are also a very important energy source. In vertebrates at least half of the oxidative energy used by the liver, kidneys, heart muscle, and resting skeletal muscle is derived from the oxidation of fatty acids; in fasting or hibernating animals or in migrating birds, fat is virtually the sole source of energy.

Neutral fats or triglycerides, the major components of storage fats in plant and animal cells, consist of the alcohol glycerol linked to three molecules of fatty acids. Before a molecule of neutral fat can be metabolized, it must be hydrolyzed to its component parts. Hydrolysis [19] is effected by intracellular enzymes or gut enzymes, and forms phase I of fat catabolism. Letters x , y , and z represent the number of $-\text{CH}_2-$ groups in the fatty acid molecules.

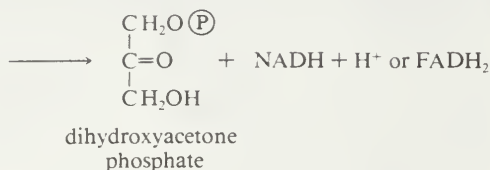
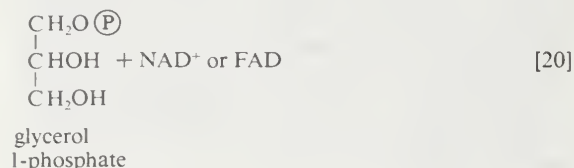
Tri-glycerides



As is apparent from [19], the three molecules of fatty acid released from the triglyceride need not be identical; a fatty acid usually contains 16 or 18 carbon atoms but may also be unsaturated; that is, containing one or more double bonds ($-\text{CH}=\text{CH}-$). Only the fate of saturated fatty acids, of the type $\text{CH}_3(\text{CH}_2)_n\text{COOH}$ (n most commonly is an even number), is dealt with here.

Fate of glycerol. It requires but two reactions to channel glycerol into a catabolic pathway (see Figure 2). In a reaction catalyzed by glycerolkinase, ATP is used to phosphorylate glycerol; the products are glycerol 1-phosphate and ADP. Glycerol 1-phosphate is then oxidized to dihydroxyacetone phosphate [20], an intermediate of glycolysis. The reaction is catalyzed by either a soluble (cytoplasmic) enzyme, glycerolphosphate dehydrogenase, or a similar enzyme present in the mitochondria. In addition to their different locations, the two dehydrogenase enzymes differ in that a different coenzyme accepts the electrons removed

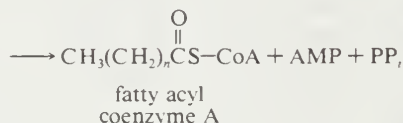
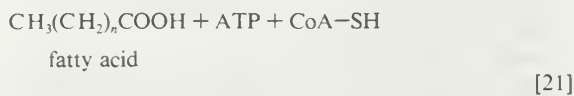
from glycerol 1-phosphate. In the case of the cytoplasmic enzyme, NAD⁺ accepts the electrons (and is reduced to NADH + H⁺); in the case of the mitochondrial enzyme, flavin adenine dinucleotide (FAD) accepts the electrons (and is reduced to FADH₂).



Fate of fatty acids. As with sugars, the release of energy from fatty acids necessitates an initial investment of ATP. A problem unique to fats is a consequence of the low solubility in water of most fatty acids. Their catabolism requires mechanisms that fragment them in a controlled and stepwise manner. The mechanism involves a coenzyme for the transfer of an acyl group

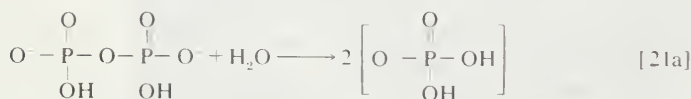
Coenzyme A and the fragmentation of fats

(e.g., CH₃C=O), namely, coenzyme A. The functional portion of this complex molecule is the sulfhydryl (-SH) group at one end. The coenzyme is often identified as CoA-SH (see step [21]). The organized and stepwise degradation of fatty acids linked to coenzyme A is ensured because the necessary enzymes are sequestered in particulate structures. In microorganisms these enzymes are associated with cell membranes, and in higher organisms with mitochondria.

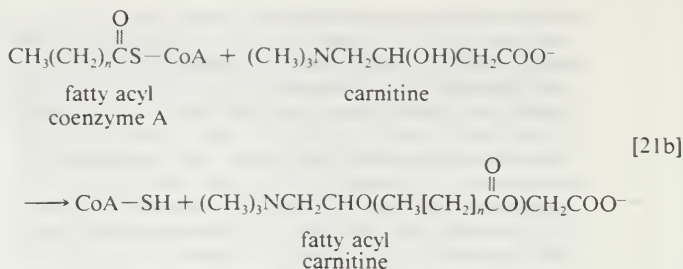


Fatty acids are linked to coenzyme A (CoA-SH) in one of two main ways. In higher organisms, enzymes in the cytoplasm called thiokinases catalyze the linkage of fatty acids with CoA-SH to form a compound that can be called a fatty acyl coenzyme A [21]. This step requires ATP, which is split to AMP and inorganic pyrophosphate (PP_i) in the process.

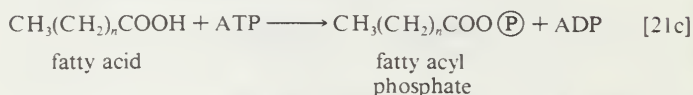
In this series of reactions, *n* indicates the number of hydrocarbon units (-CH₂-) in the molecule. Because most tissues contain highly active pyrophosphatase enzymes [21a], which catalyze the virtually irreversible hydrolysis of inorganic pyrophosphate (PP_i) to two molecules of inorganic phosphate (P_i), reaction [21] proceeds overwhelmingly to completion; i.e., from left to right.



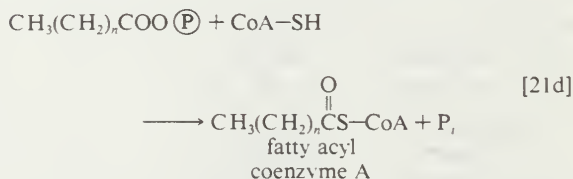
Although fatty acids are activated in this way, the acyl coenzyme A derivatives that are formed must be transported to the enzyme complex that effects their oxidation. Activation occurs in the cytoplasm, but, in animal cells, oxidation takes place in the mitochondria. The transfer of fatty acyl coenzyme A across the mitochondrial membrane is effected by the enzyme carnitine, a nitrogen-containing small hydroxy acid of the formula (CH₃)₃NCH₂CH(OH)CH₂COO⁻. The -OH group within



the carnitine molecule accepts the acyl group of fatty acyl coenzyme A, forming acyl carnitine, which can cross the inner membrane of the mitochondrion and there return the acyl group to coenzyme A.



These reactions are catalyzed by the enzyme carnitine acyl transferase. Defects in this enzyme or in the carnitine carrier are inborn errors of metabolism. In obligate anaerobic bacteria the linkage of fatty acids to coenzyme A may require the formation of a fatty acyl phosphate, i.e., the phosphorylation of the fatty acid using ATP; ADP is also a product [21c]. The fatty acyl moiety [CH₃(CH₂)_nCOO⁻] is then transferred to coenzyme A [21d], forming a fatty acyl coenzyme A compound and P_i. In either case, it is the fatty acyl coenzyme A molecules that are fragmented in the sequence of events summarized in Figure 5.



Initially (step [22], Figure 5) two hydrogen atoms are lost from the fatty acyl coenzyme A, resulting in the formation of an unsaturated fatty acyl coenzyme A (i.e., with a double bond, -CH=CH-) between the α- and β-carbons of the acyl moiety. (The α-carbon is the one closest to

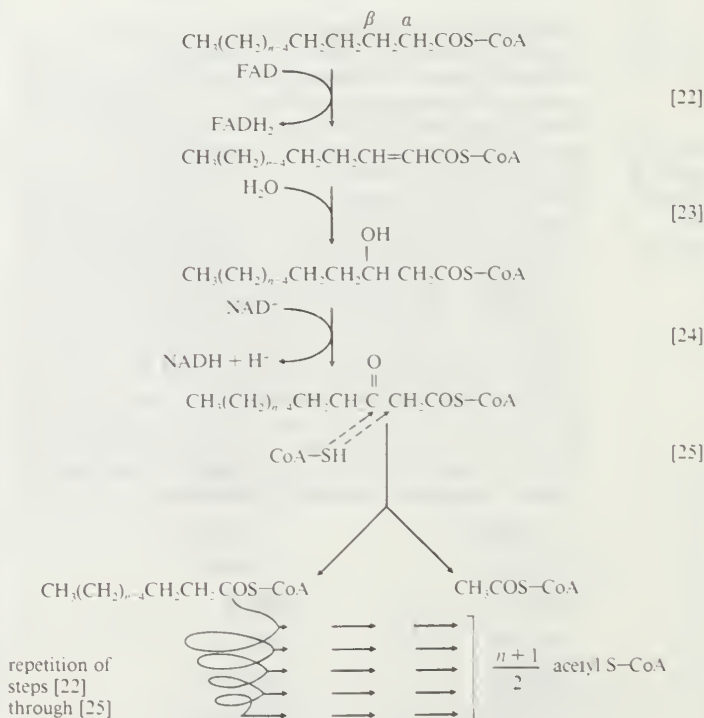


Figure 5: Fragmentation of acyl coenzyme A with an even number of carbon atoms (see text).

the carboxyl [-COOH] group of a fatty acid; the next closest is the β -, and so on to the end of the hydrocarbon chain.) The hydrogen atoms are accepted by the coenzyme FAD (flavin adenine dinucleotide), which is reduced to FADH₂. The product of step [22], α,β -unsaturated fatty acyl coenzyme A, is enzymatically hydrated [23]; *i.e.*, water is added across the double bond. The product, called a β -hydroxyacyl coenzyme A, can again be oxidized in an enzyme-catalyzed reaction [24]; the electrons removed are accepted by NAD⁺. The product is called a β -ketoacyl coenzyme A.

The next enzymatic step [25] enables the energy invested in step [21] to be conserved. The β -ketoacyl coenzyme A that is the product of reaction [24] is split, not by water but by coenzyme A. The process, called thiolysis (as distinct from hydrolysis), yields the two-carbon fragment acetyl coenzyme A and a fatty acyl coenzyme A having two fewer carbon atoms than the molecule that underwent reaction [22]; otherwise the two are similar.

The shortened fatty acyl coenzyme A molecule now undergoes the sequence of reactions again, beginning with the dehydrogenation step [22], and another two-carbon fragment is removed as acetyl coenzyme A. With each passage through the process of fatty acid oxidation, the fatty acid loses a two-carbon fragment as acetyl coenzyme A and two pairs of hydrogen atoms to specific acceptors. The 16-carbon fatty acid, palmitic acid, for example, undergoes a total of seven such cycles, yielding eight molecules of acetyl coenzyme A and 14 pairs of hydrogen atoms, seven of which appear in the form of FADH₂ and seven in the form of NADH + H⁺. The reduced coenzymes, FADH₂ and reduced NAD⁺, are reoxidized when the electrons pass through the electron transport chain, with concomitant formation of ATP (see below *Biological energy transduction*). In anaerobes, organic molecules and not oxygen are electron acceptors; thus the yield of ATP is reduced. In all organisms, however, the acetyl coenzyme A formed from the breakdown of fatty acids joins that arising from the catabolism of carbohydrates (see below *The oxidation of pyruvate*) and many amino acids (see below *The catabolism of proteins*); Figure 2 shows the interrelationships.

Fatty acids with an odd number of carbon atoms are relatively rare in nature but may arise during microbial fermentations or through the oxidation of amino acids such as valine and isoleucine. They may be fragmented through repeated cycles of steps [22] to [25] until the final five-carbon acyl coenzyme A is split into acetyl coenzyme A and propionyl coenzyme A, which has three carbon atoms. In many bacteria this propionyl coenzyme A can be transformed either to acetyl coenzyme A and carbon dioxide or to pyruvate. In other microorganisms and in animals propionyl coenzyme A has a different fate: carbon dioxide is added to propionyl coenzyme A in a reaction requiring ATP. The product, methylmalonyl coenzyme A, has four carbon atoms; the molecule undergoes a rearrangement, forming succinyl coenzyme A, which is an intermediate of the TCA cycle.

THE CATABOLISM OF PROTEINS

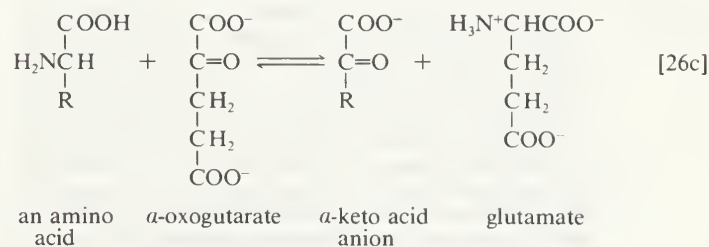
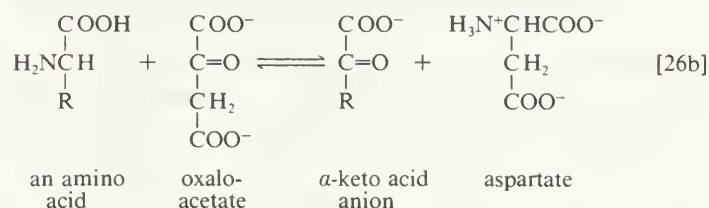
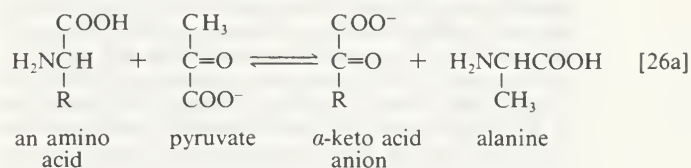
Role of amino acids

The amino acids derived from proteins function primarily as the precursors, or building blocks, for the cell's own proteins and (unlike lipids and carbohydrates) are not primarily a source of energy. Many microorganisms, on the other hand, can grow by using amino acids as the sole carbon and nitrogen source. Under these conditions these microorganisms derive from the amino acids all of their required energy and all of the precursors of the macromolecules that comprise the components of their cells. Moreover, it has been calculated that a man of average weight (70 kilograms, or 154 pounds) turns over about 0.4 kilogram of protein per day. About 0.1 kilogram is degraded and replaced by dietary amino acids; the remaining 0.3 kilogram is recycled as part of the dynamic state of cell constituents. The cells of plants contain and metabolize many amino acids in addition to the 20 or so that are normally found in proteins. A complete discussion of these special pathways is outside the scope of this article, however.

Before proteins can enter cells, the bonds linking adjacent amino acids (peptide bonds) must be hydrolyzed; this process releases the amino acids constituting the protein. The utilization of dietary proteins thus requires the operation of extracellular digestive enzymes; *i.e.*, enzymes outside the cell. Many microorganisms secrete such enzymes into the nutrient media in which they are growing; animals secrete them into the gut. The turnover of proteins within cells, on the other hand, requires the functioning of intracellular enzymes that catalyze the splitting of the peptide bonds linking adjacent amino acids; little is known about the mechanism involved.

Amino acids may be described by the general formula RCH(NH₂)COOH, or RCH(NH₃⁺)COO⁻, in which R represents a specific chemical moiety. The catabolic fate of amino acids involves (1) removal of nitrogen, (2) disposal of nitrogen, and (3) oxidation of the remaining carbon skeleton.

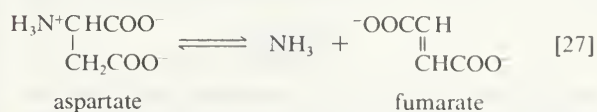
Removal of nitrogen. The removal of the amino group (-NH₂) generally constitutes the first stage in amino-acid catabolism. The amino group usually is initially transferred to the anion of one of three different α -keto acids (*i.e.*, of the general structure RCOCOO⁻): pyruvate, which is an intermediate of carbohydrate fragmentation; or oxaloacetate or α -oxoglutarate, both intermediates of the TCA cycle. The products are alanine, aspartate, and glutamate (reactions [26a, b, and c]).



Since the effect of these reactions is to produce n amino acids and n keto acids from n different amino acids and n different keto acids, no net reduction in the nitrogen content of the system has yet been achieved. The elimination of nitrogen occurs in a variety of ways.

Elimination of nitrogen

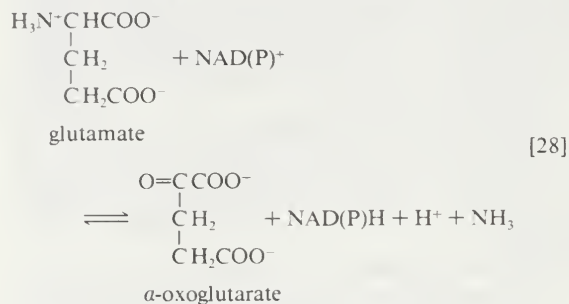
In many microorganisms, ammonia (NH₃) can be removed from aspartate via a reaction catalyzed by aspartase [27]; the other product, fumarate, is an intermediate of the TCA cycle.



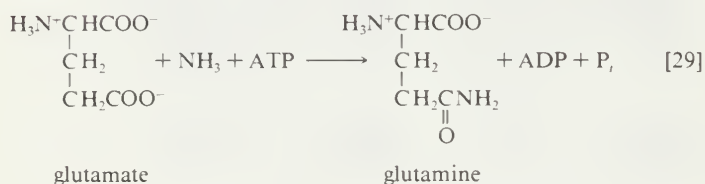
A quantitatively more important route is that catalyzed by glutamate dehydrogenase, in which the glutamate formed in [26c] is oxidized to α -oxoglutarate, another TCA cycle intermediate [28]. Either NADP⁺ or both NADP⁺ and NAD⁺ may serve as the hydrogen or electron acceptor, depending on the organism; and some organisms synthesize two enzymes, one of which prefers NADP⁺ and the

other NAD^+ . In reaction [28], NAD(P)^+ is used to indicate that either NAD^+ , NADP^+ , or both may serve as the electron acceptor.

The occurrence of the transfer reactions [26] and either step [27] or, more importantly, step [28] allows the channeling of many amino acids into a common pathway by which nitrogen can be eliminated as ammonia.

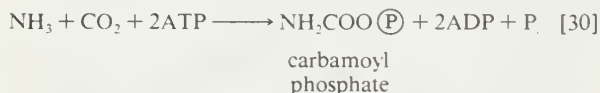


Disposal of nitrogen. In animals that excrete ammonia as the main nitrogenous waste product (*e.g.*, some marine invertebrates, crustaceans), it is derived from nitrogen transfer reactions [26] and oxidation via glutamate dehydrogenase [28] as described above for microorganisms. Because ammonia is toxic to cells, however, it is detoxified as it forms. This process involves an enzyme-catalyzed reaction between ammonia and a molecule of glutamate; ATP provides the energy for the reaction, which results in the formation of glutamine. ADP, and inorganic phosphate [29]. This reaction [29] is catalyzed by glutamine synthetase, which is subject to a variety of metabolic controls. The glutamine thus formed gives up the amide nitrogen in the kidney tubules. As a result glutamate is formed once again, and ammonia is released into the urine.

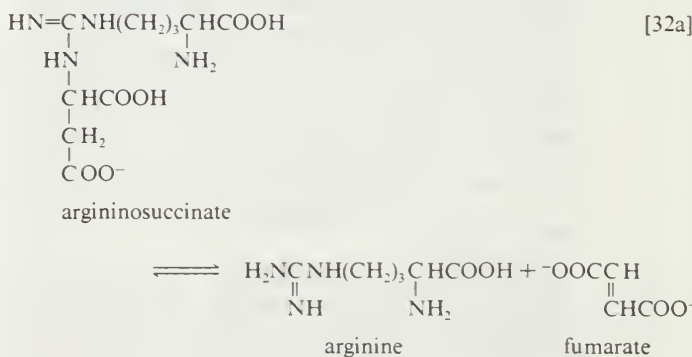
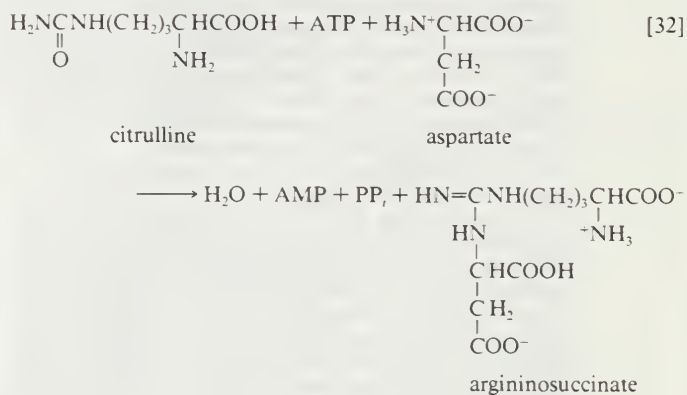
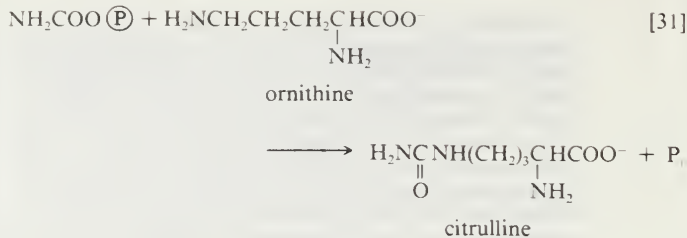


Uric acid In terrestrial reptiles and birds, uric acid rather than glutamate is the compound with which nitrogen combines to form a nontoxic substance for transfer to the kidney tubules. Uric acid is formed by a complex pathway that begins with ribose 5-phosphate and during which a so-called purine skeleton (see Figure 11) is formed; in the course of this process, nitrogen atoms from glutamine and the amino acids aspartic acid and glycine are incorporated into the skeleton. These nitrogen donors are derived from other amino acids via amino group transfer [26] and the reaction catalyzed by glutamine synthetase [29].

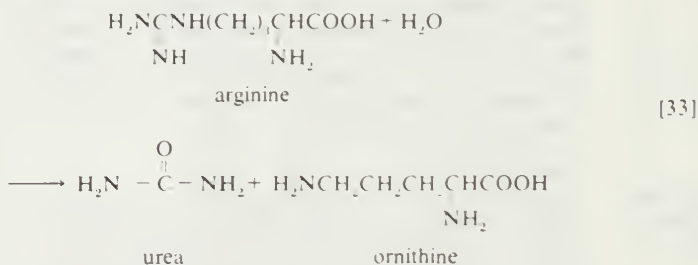
In most fishes, amphibians, and mammals, nitrogen is detoxified in the liver and excreted as urea, a readily soluble and harmless product. The sequence leading to the formation of urea, commonly called the urea cycle, is summarized as follows: Ammonia, formed from glutamate and NAD^+ in the liver mitochondria (reaction [28]), reacts with carbon dioxide and ATP to form carbamoyl phosphate, ADP, and inorganic phosphate, as shown in reaction [30].



The reaction is catalyzed by carbamoyl phosphate synthetase. The carbamoyl moiety of carbamoyl phosphate ($\text{NH}_2\text{CO}-$) is transferred to ornithine, an amino acid, in a reaction catalyzed by ornithine transcarbamoylase; the products are citrulline and inorganic phosphate [31]. Citrulline and aspartate formed from amino acids via step [26b] react to form argininosuccinate [32]; argininosuccinic acid synthetase catalyzes the reaction. Argininosuccinate splits into fumarate and arginine during a reaction



catalyzed by argininosuccinase [32a]. In the final step of the urea cycle, arginine, in a reaction catalyzed by arginase, is hydrolyzed [33]. Urea and ornithine are the products; ornithine thus is available to initiate another cycle beginning at step [31].

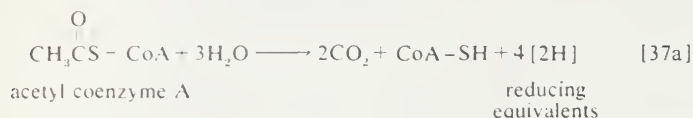


Oxidation of the carbon skeleton. As indicated in Figure 2, the carbon skeletons of amino acids (*i.e.*, the portion of the molecule remaining after the removal of nitrogen) are fragmented to form only a few end products; all of them are intermediates of either glycolysis or the TCA cycle. The number and complexity of the catabolic steps by which each amino acid arrives at its catabolic end point reflects the chemical complexity of that amino acid. Thus, in the case of alanine, only the amino group must be removed to yield pyruvate; the amino acid threonine, on the other hand, must be transformed successively to the amino acids glycine and serine before pyruvate is formed. The fragmentation of leucine to acetyl coenzyme A involves seven steps; that of tryptophan to the same end product requires 11. (A detailed discussion of the events

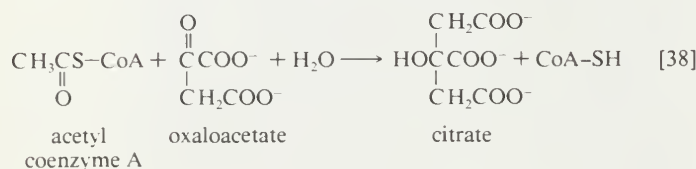
End
products
as inter-
mediates

Citrulline

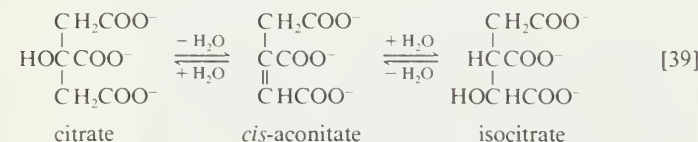
catalyzed steps that effects the total combustion of the acetyl moiety of the coenzyme represents the terminal oxidative pathway for virtually all food materials. The balance of the overall reaction of the TCA cycle [37a] is that three molecules of water react with acetyl coenzyme A to form carbon dioxide, coenzyme A, and reducing equivalents. The oxidation by oxygen of the reducing equivalents is accompanied by the conservation (as ATP) of most of the energy of the food ingested by aerobic organisms.



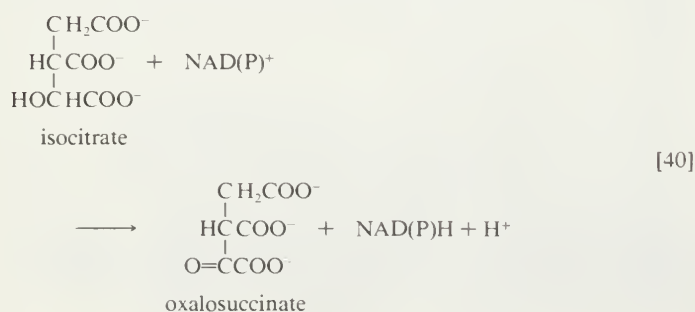
The relative complexity and number of chemical events that constitute the TCA cycle, and their location as components of spatially determined structures such as cell membranes in microorganisms and mitochondria in plants and higher animals, reflect the problems involved chemically in "dismembering" a compound having only two carbon atoms and releasing in a controlled and stepwise manner the reducing equivalents ultimately to be passed to oxygen. These problems have been overcome by the simple but effective device of initially combining the two-carbon compound with a four-carbon acceptor; it is much less difficult chemically to dismember and oxidize a compound having six carbon atoms.



In the TCA cycle, acetyl coenzyme A initially reacts with oxaloacetate to yield citrate and to liberate coenzyme A. This reaction [38] is catalyzed by citrate synthase. (As mentioned above, many of the compounds in living cells that take part in metabolic pathways exist as charged moieties, or anions, and are named as such.) Citrate undergoes isomerization (*i.e.*, a rearrangement of certain atoms comprising the molecule) to form isocitrate [39]. The reaction involves first the removal of the elements of water from citrate to form *cis*-aconitate, and then the re-addition of water to *cis*-aconitate in such a way that isocitrate is formed. It is probable that all three reactants—citrate, *cis*-aconitate, and isocitrate—remain closely associated with aconitase, the enzyme that catalyzes the isomerization process, and that most of the *cis*-aconitate is not released from the enzyme surface but is immediately converted to isocitrate.

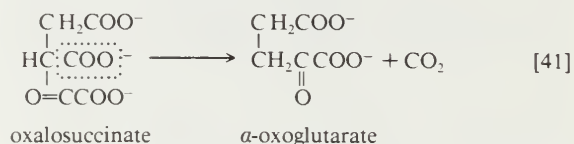


Isocitrate is oxidized—*i.e.*, hydrogen is removed—to form oxalosuccinate; the two hydrogen atoms are usually transferred to NAD^+ , thus forming reduced NAD^+ [40]; in some microorganisms, and during the biosynthesis of glutamate in the cytoplasm of animal cells, however, the

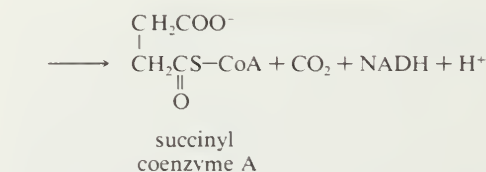
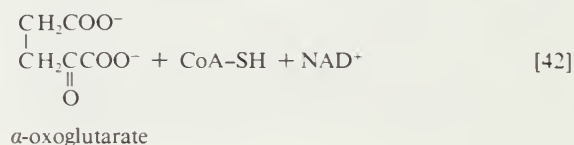


hydrogen atoms may also be accepted by NADP^+ . Thus the enzyme controlling this reaction, isocitrate dehydrogenase, differs in specificity for the coenzymes; various forms occur not only in different organisms but even within the same cell. In [40] NAD(P)^+ indicates that either NAD^+ or NADP^+ can act as a hydrogen acceptor.

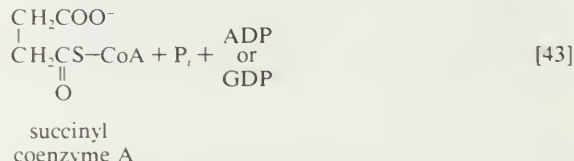
The position of the carboxylate ($-\text{COO}^-$) that is "sandwiched" in the middle of the oxalosuccinate molecule renders it very unstable; as a result, the carbon of this group is lost as carbon dioxide (note the dotted rectangle) in a reaction [41] that can occur spontaneously but may be further accelerated by an enzyme.



The five-carbon product of reaction [41], α -oxoglutarate, has chemical properties similar to pyruvate (free-acid forms of both are so-called α -oxoacids), and the chemical events involved in the oxidation of α -oxoglutarate are analogous to those already described for the oxidation of pyruvate (see reaction [37]). Reaction [42] is effected by a multi-enzyme complex; TPP, lipS₂ (6,8-dithio-*n*-octanoate), and coenzyme A are required as coenzymes. The products are carbon dioxide and succinyl coenzyme A. As was noted with reaction [37], this oxidation of α -oxoglutarate results in the reduction of lipS₂, which must be reoxidized. This is done by transfer of reducing equivalents to FAD and thence to NAD^+ . The resultant $\text{NADH} + \text{H}^+$ is reoxidized by the passage of the electrons, ultimately, to oxygen, via the electron transport chain.



Unlike the acetyl coenzyme A produced from pyruvate in reaction [37], succinyl coenzyme A undergoes a phosphorylation reaction—*i.e.*, transfer of the succinyl moiety from coenzyme A to inorganic phosphate. The succinyl phosphate thus formed is not released from the enzyme surface; an unstable, high-energy compound called an acid anhydride, it transfers a high-energy phosphate to ADP, directly or via guanosine diphosphate (GDP) [43].



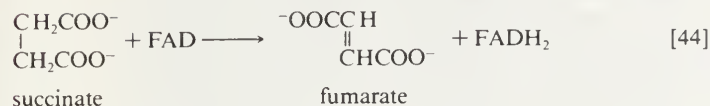
Phosphorylation

If guanosine triphosphate (GTP) forms, ATP can readily arise from it in an exchange involving ADP [43a]:

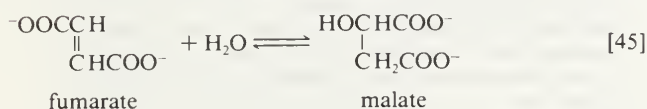


The remainder of the reactions of the TCA cycle serve to regenerate the initial four-carbon acceptor of acetyl coenzyme A (oxaloacetate) from succinate, the process

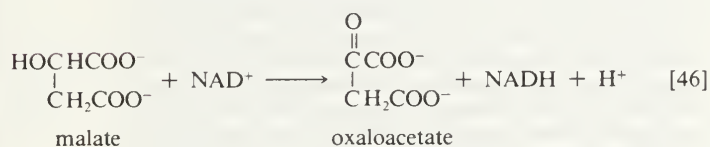
requiring in effect the oxidation of a methylene group ($-\text{CH}_2-$) to a carbonyl group ($-\text{CO}-$), with concomitant release of $2 \times [2\text{H}]$ reducing equivalents. It is therefore similar to, and is effected in like manner to, the oxidation of fatty acids (steps [22–24]; see Figure 5). As is the case with fatty acids, hydrogen atoms or electrons are initially removed from the succinate formed in [43] and are accepted by FAD; the reaction, catalyzed by succinate dehydrogenase [44], results in the formation of fumarate and reduced FAD.



The elements of water are added across the double bond ($-\text{CH}=\text{CH}-$) of fumarate in a reaction catalyzed by fumarase [45]; this type of reaction also occurred in step [39] of the cycle. The product of reaction [45] is malate.



Malate can be oxidized to oxaloacetate by removal of two hydrogen atoms, which are accepted by NAD^+ . This type of reaction, catalyzed by malate dehydrogenase in reaction [46], also occurred in step [40] of the cycle. The formation of oxaloacetate completes the TCA cycle, which can now begin again with the formation of citrate [38].



Energy conserved during the TCA cycle The loss of the two molecules of carbon dioxide in steps [41] and [42] does not yield biologically useful energy. The substrate-linked formation of ATP accompanies step [43], in which one molecule of ATP, is formed during each turn of the cycle. The hydrogen ions and electrons that result from steps [40], [42], [44], and [46] are passed down the chain of respiratory carriers to oxygen, with the concomitant formation of three molecules of ATP, per $[2\text{H}]$, as $\text{NADH} + \text{H}^+$ (see below). Similarly, the oxidation of the reduced FAD formed in [44] results in the formation of two ATP. Each turn of the cycle thus leads to the production of a total of 12 ATP. It will be recalled that the anaerobic fragmentation of glucose to two molecules of pyruvate yielded two ATP; the aerobic oxidation via the TCA cycle of two molecules of pyruvate thus makes available to the cell at least 15 times more ATP per molecule of glucose catabolized than is produced anaerobically. If, in addition, the $2 \times [\text{NADH} + \text{H}^+]$ generated per glucose in reaction [6] are passed on to oxygen, a further six ATP are generated. The advantage to living organisms is to be able to respire rather than merely to ferment.

BIOLOGICAL ENERGY TRANSDUCTION

Adenosine triphosphate as the currency of energy exchange. When the terminal phosphate group is removed from ATP by hydrolysis, two negatively charged products are formed, ADP^{3-} and HPO_4^{2-} (a phosphate group) [47].



These products are electrically more stable than the parent molecule and do not readily recombine. The total free energy (G) of the products is much less than that of ATP; hence energy is liberated (*i.e.*, the reaction is exergonic). The amount of energy liberated under strictly defined conditions is called the standard free energy change ($\Delta G'$); this value for the hydrolysis of ATP is relatively high, at -8 kilocalories per mole. (One kilocalorie is the amount of heat required to raise the temperature of 1,000 grams of water one degree centigrade.) Conversely, the formation of ATP from ADP and inorganic phosphate (P_i) is an energy-

requiring (*i.e.*, endergonic) reaction with a standard free energy change of $+8$ kilocalories per mole.

The hydrolysis of the remaining phosphate-to-phosphate bond of ADP is also accompanied by a liberation of free energy (the standard free energy change is -6.5 kilocalories per mole); AMP hydrolysis liberates less energy (the standard free energy change is -2.2 kilocalories per mole).

The free energy of hydrolysis of a compound thus is a measure of the difference in energy content between the starting substances (reactants) and the final substances (products). Adenosine triphosphate does not have the highest standard free energy of hydrolysis of all the naturally occurring phosphates but instead occupies a position at approximately the halfway point in a series of phosphate compounds with a wide range of standard free energies of hydrolysis. Compounds such as 1,3-diphosphoglycerate or phosphoenolpyruvate (PEP), which are above ATP on the scale (see Figure 6), have large negative $\Delta G'$ values on hydrolysis and are often called high-energy phosphates: they are said to exhibit a high phosphate group transfer potential because they have a tendency to lose their phosphate groups. Compounds such as glucose 6-phosphate or fructose 6-phosphate, which are below ATP on the scale because they have smaller negative $\Delta G'$ values on hydrolysis, have a tendency to hold on to their phosphate groups and thus act as low-energy phosphate acceptors.

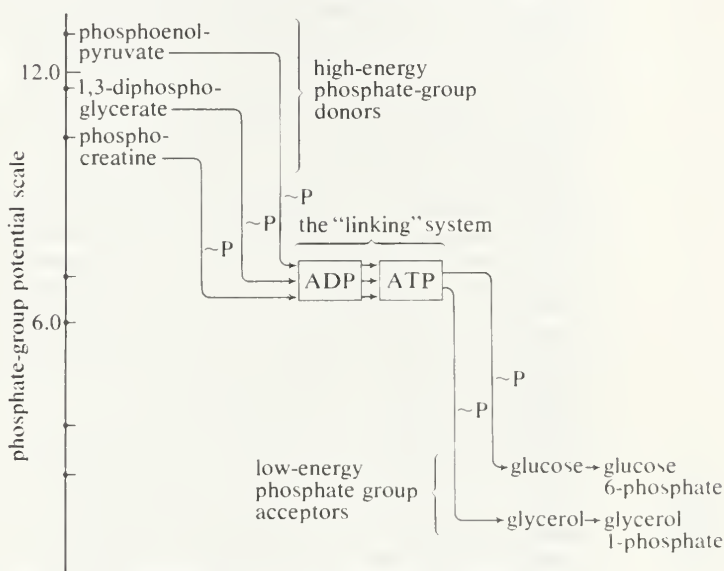
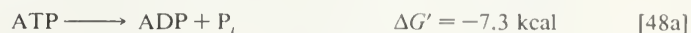
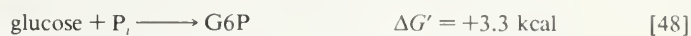


Figure 6: The transfer of phosphate groups from high-energy donors to low-energy acceptors by way of the ATP-ADP system (see text).

Both ATP and ADP act as intermediate carriers for the transfer of phosphate groups (which are more precisely called phosphoryl groups) and hence of energy, from compounds lying above ATP to those lying beneath it. Thus, in glycolysis, ADP acts as an acceptor of a phosphate group during the synthesis of ATP from PEP (see reaction [10]), and ATP functions as a donor of a phosphate group during the formation of fructose 1,6-diphosphate from fructose 6-phosphate (see reaction [3]).

The first step in glycolysis, the formation of glucose 6-phosphate (G6P), illustrates how an energetically unfavourable reaction may become feasible under intracellular conditions by coupling it to ATP.



Reaction [48] has a positive $\Delta G'$ value, indicating that the reaction tends to proceed in the reverse direction. It is therefore necessary to use the standard free energy generated by the breaking of the first phosphate bond in ATP (reaction [48a]), which is -7.3 kilocalories per mole, to move reaction [48] in the forward direction. Combining these reactions and their standard free energies gives reaction [48b] and a standard free energy value of -4 kilocalo-

ries per mole, indicating that the reaction will proceed in the forward direction. There are many intracellular reactions in which the formation of ADP or AMP from ATP provides energy for otherwise unfavourable biosyntheses. Some cellular reactions use equivalent phosphorylated analogues of ATP, for example, guanosine triphosphate (GTP) for protein synthesis.

The utilization of ATP to perform work

The function of ATP as a common intermediate of energy transfer during anabolism is further dealt with below (see *The biosynthesis of cell components*). In certain specialized cells or tissues the chemical energy of ATP is used to perform work other than the chemical work of anabolism; for example, mechanical work—such as muscular contraction, or the movement of contractile structures called cilia and flagella, which are responsible for the motility of many small organisms. The performance of osmotic work also requires ATP; *e.g.*, the transport of ions or metabolites through membranes against a concentration gradient, a process that is basically responsible for many physiological functions, including nerve conduction, the secretion of hydrochloric acid in the stomach, and the removal of water from the kidneys.

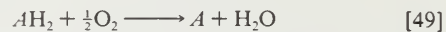
Energy conservation. The amount of ATP in a cell is limited, and it must be replaced continually to maintain repair and growth. This is achieved by using the energy liberated during the oxidative stages of catabolism to synthesize ATP from ADP and phosphate. The synthesis of ATP linked to catabolism occurs by two distinct mechanisms: substrate-level phosphorylation and oxidative, or respiratory-chain, phosphorylation. Oxidative phosphorylation is the major method of energy conservation under aerobic conditions in all nonphotosynthetic cells.

Substrate-level phosphorylation. In substrate-level phosphorylation a phosphoryl group is transferred from an energy-rich donor (*e.g.*, 1,3-diphosphoglycerate) to ADP to yield a molecule of ATP. This type of ATP synthesis (see reactions [7], [10], and [43]) does not require molecular oxygen (O_2), although it is frequently, but not always, preceded by an oxidation (*i.e.*, dehydrogenation) reaction. Substrate-level phosphorylation is the major method of energy conservation in oxygen-depleted tissues and during fermentative growth of microorganisms.

Oxidative, or respiratory-chain, phosphorylation. In oxidative phosphorylation the oxidation of catabolic intermediates by molecular oxygen occurs via a highly ordered series of substances that act as hydrogen and electron carriers. They constitute the electron transfer system, or respiratory chain. In most animals, plants, and fungi, the electron transfer system is fixed in the membranes of mitochondria; in bacteria (which have no mitochondria) this system is incorporated into the plasma membrane. Sufficient free energy is released to allow the synthesis of ATP by a process described below. First, however, it is necessary to consider the nature of the respiratory chain.

The electron transfer system

Four types of hydrogen or electron carriers are known to participate in the respiratory chain, in which they serve to transfer two reducing equivalents (2H) from reduced substrate (AH_2) to molecular oxygen (see reaction [49]); the products are the oxidized substrate (A) and water (H_2O).



The carriers are NAD^+ and, less frequently, $NADP^+$; the flavoproteins FAD and FMN (flavin mononucleotide); ubiquinone (or coenzyme Q); and several types of cytochromes. Each carrier has an oxidized and reduced form (*e.g.*, FAD and $FADH_2$, respectively), the two forms constituting an oxidation-reduction, or redox, couple. Within the respiratory chain each redox couple undergoes cyclic oxidation-reduction—*i.e.*, the oxidized component of the couple accepts reducing equivalents from either a substrate or a reduced carrier preceding it in the series and in turn donates these reducing equivalents to the next oxidized carrier in the sequence. Reducing equivalents are thus transferred from substrates to molecular oxygen by a number of sequential redox reactions.

Most oxidizable catabolic intermediates initially undergo a dehydrogenation reaction, during which a dehydrogenase enzyme transfers the equivalent of a hydride ion ($H^+ + 2e^-$, with e^- representing an electron) to its coenzyme, either NAD^+ or $NADP^+$. The reduced NAD^+ (or $NADP^+$) thus produced (usually written as $NADH + H^+$ or $NADPH + H^+$) diffuses to the membrane-bound respiratory chain to be oxidized by an enzyme known as NADH dehydrogenase; the enzyme has as its coenzyme FMN. There is no corresponding NADPH dehydrogenase in mammalian mitochondria; instead, the reducing equivalents of $NADPH + H^+$ are transferred to NAD^+ in a reaction catalyzed by a transhydrogenase enzyme, with the products being reduced $NADH + H^+$ and $NADP^+$. A few substrates (*e.g.*, acyl coenzyme A and succinate; see reactions [22] and [44]) bypass this reaction and instead undergo immediate dehydrogenation by specific membrane-bound dehydrogenase enzymes. During the reaction, the coenzyme FAD accepts two hydrogen atoms and two electrons ($2H + 2e^-$). The reduced flavoproteins (*i.e.*, $FMNH_2$ and $FADH_2$) donate their two hydrogen atoms to the lipid carrier ubiquinone, which is thus reduced.

The fourth type of carrier, the cytochromes, consists of hemoproteins—*i.e.*, proteins with a nonprotein component, or prosthetic group, called heme (or a derivative of heme), which is an iron-containing pigment molecule. The iron atom in the prosthetic group is able to carry one electron and oscillates between the oxidized, or ferric (Fe^{3+}), and the reduced, or ferrous (Fe^{2+}), forms. The five cytochromes present in the mammalian respiratory chain, designated cytochromes b , c_1 , c , a , and a_3 , act in sequence between ubiquinone and molecular oxygen. The terminal cytochrome of this sequence (a_3 , also known as

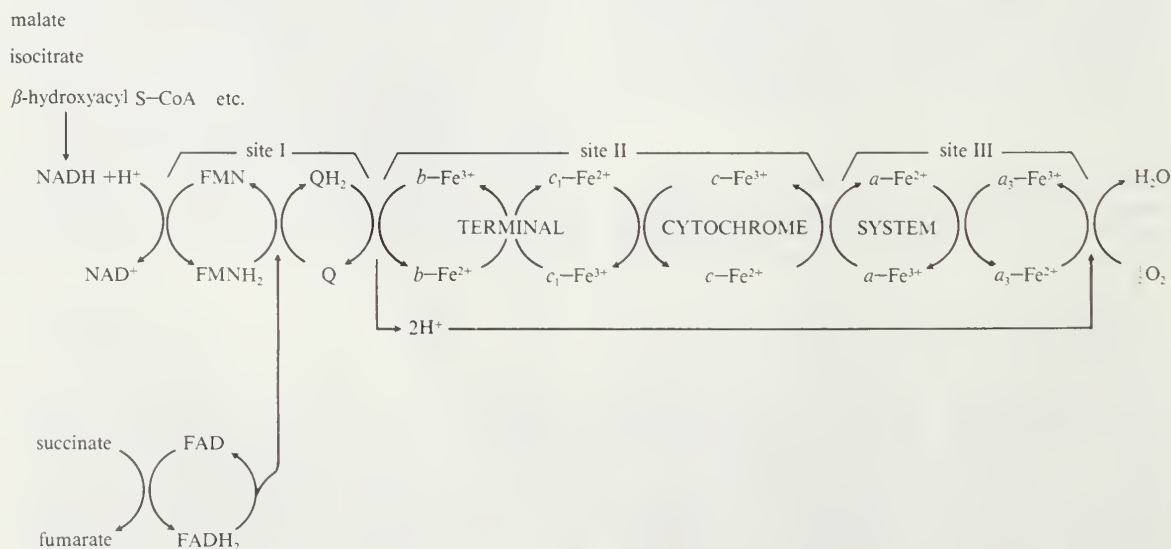


Figure 7: The respiratory chain (see text).

Sequence
of carriers

cytochrome oxidase) is able to donate electrons to oxygen rather than to another electron carrier; a_3 is also the site of action of two substances that inhibit the respiratory chain, potassium cyanide and carbon monoxide. Special Fe-S complexes play a role in the activity of NADH dehydrogenase and succinate dehydrogenase. The sequence of carriers, from substrates to oxygen, is shown schematically in Figure 7.

In each redox couple the reduced form has a tendency to lose reducing equivalents (*i.e.*, to act as an electron or hydrogen donor); similarly, the oxidized form has a tendency to gain reducing equivalents (*i.e.*, to act as an electron or hydrogen acceptor). The oxidation-reduction characteristics of each couple can be determined experimentally under well-defined, standard conditions. The value thus obtained is the standard oxidation-reduction (redox) potential (E_0). Values for respiratory chain carriers range from $E_0 = -320$ millivolts (one millivolt = 0.001 volt) for NAD^+ /reduced NAD^+ to $E_0 = +820$ millivolts for $\frac{1}{2}\text{O}_2/\text{H}_2\text{O}$; the values for intermediate carriers lie between. Reduced NAD^+ is the most electronegative carrier, oxygen the most electropositive acceptor. During respiration reducing equivalents undergo stepwise transfer from the reduced form of the most electronegative carrier (reduced NAD^+) to the oxidized form of the most electropositive couple (oxygen). Each step is accompanied by a decline in standard free energy ($\Delta G'$) proportional to the difference in the standard redox potentials (ΔE_0) of the two carriers involved.

Oxidative
phosphorylation

Overall oxidation of reduced NAD^+ by oxygen ($\Delta E_0 = +1,140$ millivolts) is accompanied by the liberation of free energy ($\Delta G' = -52.4$ kilocalories per mole); in theory this energy is sufficient to allow the synthesis of six or seven molecules of ATP. In the cell, however, this synthesis of ATP, called oxidative phosphorylation, proceeds with an efficiency of about 46 percent; thus only three molecules of ATP are produced per atom of oxygen consumed—this being the so-called P/O, P/O, or ADP/O ratio. The energy that is not conserved as ATP is lost as heat. The oxidation of succinate by molecular oxygen ($\Delta E_0 = +790$ millivolts), which is accompanied by a smaller liberation of free energy ($\Delta G' = -36.5$ kilocalories per mole), yields only two molecules of ATP per atom of oxygen consumed (P/O = 2).

In order to understand the mechanism by which the energy released during respiration is conserved as ATP, it is necessary to appreciate the structural features of mitochondria. These are organelles in animal and plant cells in which oxidative phosphorylation takes place. There are many mitochondria in animal tissues; for example, in heart and skeletal muscle, which require large amounts of energy for mechanical work, in the pancreas, where there is biosynthesis, and in the kidney, where the process of excretion begins. Mitochondria have an outer membrane, which allows the passage of most small molecules and ions, and a highly folded inner membrane (crista), which does not even allow the passage of small ions and so maintains a closed space within the cell. The electron-transferring molecules of the respiratory chain and the enzymes responsible for ATP synthesis are located in and on this inner membrane, while the space inside (matrix) contains the enzymes of the TCA cycle (reactions [34] to [46]; see also CELLS). The enzyme systems primarily responsible for the release and subsequent oxidation of reducing equivalents are thus closely related so that the reduced coenzymes formed during catabolism ($\text{NADH} + \text{H}^+$ and FADH_2) are available as substrates for respiration. The movement of most charged metabolites into the matrix space is mediated by special carrier proteins in the crista that catalyze exchange-diffusion (*i.e.*, a one-for-one exchange). The oxidative phosphorylation systems of bacteria are similar in principle but show a greater diversity in the composition of their respiratory carriers.

The mechanism of ATP synthesis appears to be as follows. During the transfer of hydrogen atoms from FMNH₂ or FADH_2 to oxygen (Figure 7), protons (H^+ ions) are pumped across the crista from the inside of the mitochondrion to the outside. Thus, respiration generates an electrical potential (and in mitochondria a small pH

gradient) across the membrane corresponding to 200 to 300 millivolts, and the chemical energy in the substrate is converted into electrical energy. Attached to the crista is a complex enzyme (ATP synthetase) that binds ATP, ADP, and P_i . It has nine polypeptide chain subunits of five different kinds in a cluster and a unit of at least three more membrane proteins composing the attachment point of ADP and P_i . This complex forms a specific proton pore in the membrane. When ADP and P_i are bound to ATP synthetase, the excess of protons (H^+) that has formed outside of the mitochondria (an H^+ gradient) moves back into the mitochondrion through the enzyme complex. The energy released is used to convert ADP and P_i to ATP. In this process, electrical energy is converted to chemical energy, and it is the supply of ADP that limits the rate of this process. The precise mechanism by which the ATP synthetase complex converts the energy stored in the electrical H^+ gradient to the chemical bond energy in ATP is not well understood. The H^+ gradient may power other endergonic (energy-requiring) processes besides ATP synthesis, such as the movement of bacterial cells and the transport of carbon substrates or ions.

Photo-
synthesis

Photosynthesis generates ATP by a mechanism that is similar in principle, if not in detail. The organelles responsible are different from mitochondria, but they also form membrane-bounded closed sacs (thylakoids) often arranged in stacks (grana). Solar energy splits two molecules of H_2O into molecular oxygen (O_2), four protons (H^+), and four electrons.

This is the source of oxygen evolution, clearly visible as bubbles from underwater plants in bright sunshine. The process involves a chlorophyll molecule, P_{680} , that changes its redox potential from +820 millivolts (in which there is a tendency to accept electrons) to about -680 millivolts (in which there is a tendency to lose electrons) upon excitation with light and acquisition of electrons. The electrons are subsequently passed along a series of carriers (plastoquinone, cytochromes b and f, and plastocyanin), analogous to the mitochondrial respiratory chain. This process pumps protons across the membrane from the outside of the thylakoid membrane to the inside. Protons (H^+) do not move freely across the membrane although chloride ions (Cl^-) do, creating a pH gradient. An ATP synthetase enzyme similar to that of the mitochondria is present, but on the outside of the thylakoid membrane. Passage of protons (H^+) through it from inside to outside generates ATP.

Hence, a gradient of protons (H^+) across the membrane is the high-energy intermediate for forming ATP in plant photosynthesis and in the respiration of all cells capable of passing reducing equivalents (hydrogen atoms or electrons) to electron acceptors.

The biosynthesis of cell components

THE NATURE OF BIOSYNTHESIS

The stages of biosynthesis. The biosynthesis of cell components (anabolism) may be regarded as occurring in two main stages. In the first, intermediate compounds of the central routes of metabolism are diverted from further catabolism and are channeled into pathways that usually lead to the formation of the relatively small molecules that serve as the building blocks, or precursors, of macromolecules.

In the second stage of biosynthesis, the building blocks are combined to yield the macromolecules—proteins, nucleic acids, lipids, and polysaccharides—that make up the bulk of tissues and cellular components. In organisms with the appropriate genetic capability, for example, all of the amino acids can be synthesized from ammonia and intermediates of the main routes of carbohydrate fragmentation and oxidation. Such intermediates act also as precursors for the purines, the pyrimidines, and the pentose sugars that constitute DNA and for a number of types of RNA. The assembly of proteins necessitates the precise combination of specific amino acids in a highly ordered and controlled manner; this in turn involves the copying, or transcription, into RNA of specific parts of DNA (see below *The synthesis of macromolecules*). The

Requirements

first stage of biosynthesis thus requires the specificity normally required for the efficient functioning of sequences of enzyme-catalyzed reactions. The second stage also involves—directly for protein and nucleic acid synthesis, less directly for the synthesis of other macromolecules—the maintenance and expression of the biological information that specifies the identity of the cell, the tissue, and the organism.

Utilization of ATP. The two stages of biosynthesis—the formation of building blocks and their specific assembly into macromolecules—are energy-consuming processes and thus require ATP. Although the ATP is derived from catabolism, catabolism does not “drive” biosynthesis. As explained in the first section of this article, the occurrence of chemical reactions in the living cell is accompanied by a net decrease in free energy. Although biological growth and development result in the creation of ordered systems from less ordered ones and of complex systems from simpler ones, these events must occur at the expense of energy-yielding reactions. The overall coupled reactions are, on balance, still accompanied by a decrease in free energy and are thus essentially irreversible in the direction of biosynthesis. The total energy released from ATP, for example, is usually much greater than is needed for a particular biosynthetic step; thus, many of the reactions involved in biosynthesis release inorganic pyrophosphate (PP_i) rather than phosphate (P_i) from ATP, and hence yield AMP rather than ADP. Since inorganic pyrophosphate readily undergoes virtually irreversible hydrolysis to two equivalents of inorganic phosphate (see [21a]), the creation of a new bond in the product of synthesis may be accompanied by the breaking of two high-energy bonds of ATP—although, in theory, one might have sufficed.

The efficient utilization for anabolic processes of ATP and some intermediate compound formed during a catabolic reaction requires the cell to have simultaneously a milieu favourable for both ATP generation and consumption. Catabolism occurs readily only if sufficient ADP is available; hence, the concentration of ATP is low. On the other hand, biosynthesis requires a high level of ATP and consequently low levels of ADP and AMP. Suitable conditions for the simultaneous function of both processes are met in two ways. Biosynthetic reactions often take place in compartments within the cell different from those in which catabolism occurs; there is thus a physical separation of energy-requiring and energy-yielding processes. Furthermore, biosynthetic reactions are regulated independently of the mechanisms by which catabolism is controlled. Such independent control is made possible by the fact that catabolic and anabolic pathways are not identical; the pacemaker, or key, enzyme that controls the overall rate of a catabolic route usually does not play any role in the biosynthetic pathway of a compound. Similarly, the pacemaker enzymes of biosynthesis are not involved in catabolism. As discussed below (see *Regulation of metabolism: Fine control*), catabolic pathways are often regulated by the relative amounts of ATP, ADP, and AMP in the cellular compartment in which the pacemaker enzymes are located; in general, ATP inhibits and ADP (or AMP) stimulates such enzymes. In contrast, many biosynthetic routes are regulated by the concentration of the end products of particular anabolic processes, so that the cell synthesizes only as much of these building blocks as it needs.

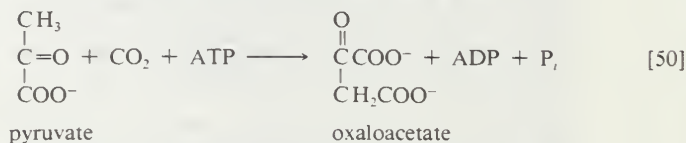
Regulation

THE SUPPLY OF BIOSYNTHETIC PRECURSORS

When higher animals consume a mixed diet, sufficient quantities of compounds for both biosynthesis and energy supply are available. Carbohydrates yield intermediates of glycolysis and of the phosphogluconate pathway, which in turn yield acetyl coenzyme A (see Figure 4); lipids yield glycolytic intermediates and acetyl coenzyme A (see Figure 2); and many amino acids form intermediates of both the TCA cycle and glycolysis. Any intermediate withdrawn for biosynthesis can thus be readily replenished by the catabolism of further nutrients. This situation does not always hold, however. Microorganisms in particular can derive all of their carbon and energy requirements by utilizing a single carbon source. The sole carbon source

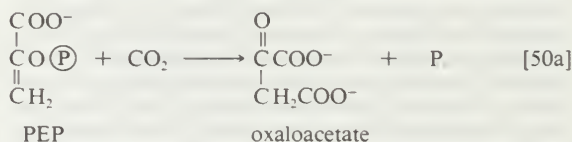
may be a substance such as a carbohydrate or a fatty acid, or an intermediate of the TCA cycle (or a substance readily converted to one). In both cases, reactions ancillary to those discussed thus far must occur before the carbon source can be utilized.

Anaplerotic routes. Although the catabolism of carbohydrates can occur via a variety of routes (see Figure 4), all give rise to pyruvate. During the catabolism of pyruvate, one carbon atom is lost as carbon dioxide and the remaining two form acetyl coenzyme A [37]; these two are involved in the TCA cycle ([41] and [42]). Because the TCA cycle is initiated by the condensation of acetyl coenzyme A with oxaloacetate, which is regenerated in each turn of the cycle, the removal of any intermediate from the cycle would cause the cycle to stop. Yet, as also indicated in Figure 4, various essential cell components are derived from *a*-oxoglutarate, succinyl coenzyme A, and oxaloacetate, so that these compounds are, in fact, removed from the cycle. Microbial growth with a carbohydrate as the sole carbon source is thus possible only if a cellular process occurs that effects the net formation of some TCA cycle intermediate from an intermediate of carbohydrate catabolism. Such a process, which replenishes the TCA cycle, has been described as an anaplerotic reaction.

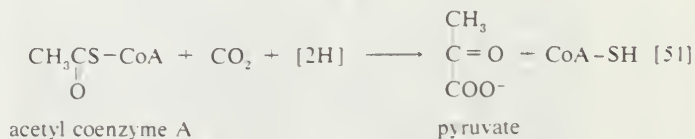


The anaplerotic function may be carried out by either of two enzymes that catalyze the fixation of carbon dioxide onto a three-carbon compound, either pyruvate [50] or phosphoenolpyruvate (PEP, [50a]) to form oxaloacetate, which has four carbon atoms. Both reactions require energy. In [50] it is supplied by the cleavage of ATP to ADP and inorganic phosphate (P_i); in [50a] it is supplied by the release of the high-energy phosphate of PEP as inorganic phosphate. Pyruvate serves as a carbon dioxide acceptor not only in many bacteria and fungi but also in the livers and kidneys of higher organisms, including man; PEP serves as the carbon dioxide acceptor in many bacteria, such as those that inhabit the gut.

Carbon dioxide acceptors

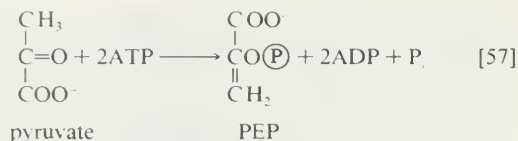


Unlike higher organisms, many bacteria and fungi can grow on acetate or compounds such as ethanol or a fatty acid that can be catabolized to acetyl coenzyme A. Under these conditions, the net formation of TCA cycle intermediates proceeds in one of two ways. In obligate anaerobic bacteria, pyruvate can be formed from acetyl coenzyme A and carbon dioxide [51]; reducing equivalents [2H] are necessary for the reaction. The pyruvate so formed can then react via either step [50] or [50a].



Reaction [51] does not occur in facultative anaerobic organisms or in strict aerobes, however. Instead, in these organisms two molecules of acetyl coenzyme A give rise to the net synthesis of a four-carbon intermediate of the TCA cycle via a route known as the glyoxylate cycle. In this route (Figure 8), the steps of the TCA cycle that lead to the loss of carbon dioxide (see [40], [41], and [42]) are bypassed. Instead of being oxidized to oxalosuccinate, as occurs in [40], isocitrate is split by isocitrate lyase [52] in a reaction similar to that of reactions [4] and [15] of carbohydrate fragmentation. The dotted line in [52] indicates the way in which isocitrate is split. The products are

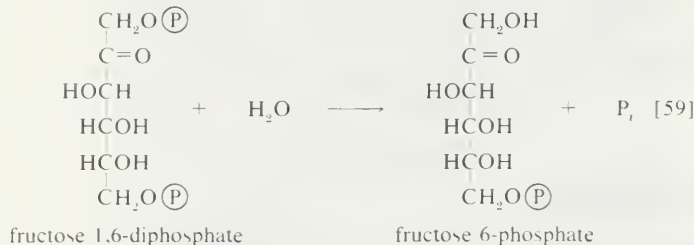
The glyoxylate cycle



The combination of steps [57] and [58] yields the same energy balance as does the direct conversion of pyruvate to PEP [56].



Phospho-fructo-kinase The second step of glycolysis bypassed in gluconeogenesis is that catalyzed by phosphofructokinase [3]. Instead, the fructose 1,6-diphosphate synthesized from dihydroxyacetonephosphate and glyceraldehyde 3-phosphate in the reaction catalyzed by aldolase is hydrolyzed, with the loss of the phosphate group linked to the first carbon atom.



The enzyme fructose diphosphatase catalyzes the reaction [59], in which the products are fructose 6-phosphate and inorganic phosphate. The fructose 6-phosphate thus formed is a precursor of mucopolysaccharides (polysaccharides with nitrogen-containing components). In addition, its conversion to glucose 6-phosphate provides the starting material for the formation of storage polysaccharides such as starch and glycogen, of monosaccharides other than glucose, of disaccharides (carbohydrates with two sugar components), and of some structural polysaccharides (*e.g.*, cellulose). The maintenance of the glucose content of

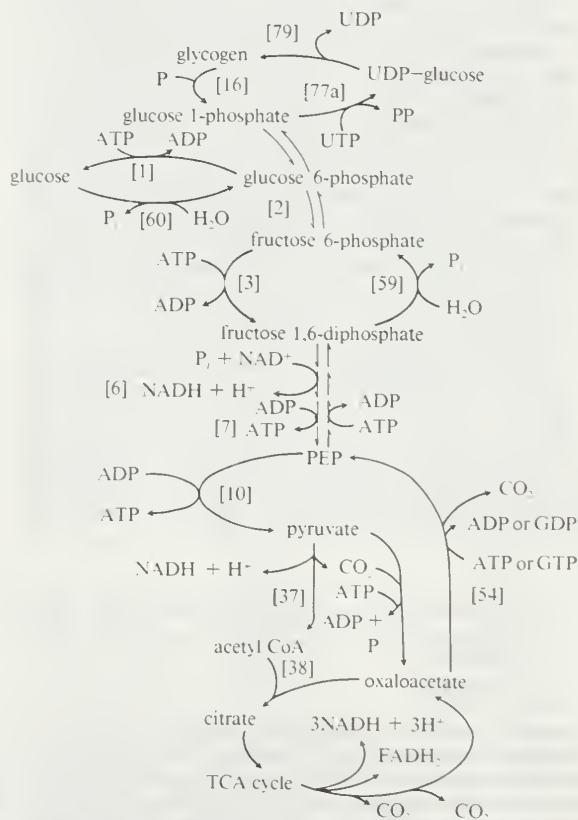
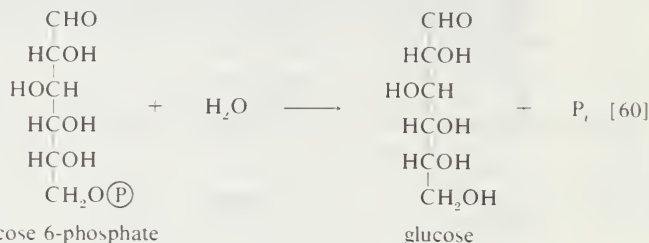


Figure 9: Catabolism and biosynthesis of glucose and glycogen. At left, reactions peculiar to catabolism. At right, reactions peculiar to anabolism. Numbers in brackets refer to reactions explained in text.

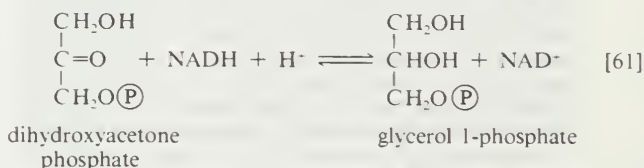
vertebrate blood requires glucose 6-phosphate to be converted to glucose. This process occurs in the kidney, in the lining of the intestine, and most importantly in the liver. The reaction does not occur by reversal of the hexokinase or glucokinase reactions that effect the formation of glucose 6-phosphate from glucose and ATP [1]; rather, glucose 6-phosphate is hydrolyzed in a reaction catalyzed by glucose 6-phosphatase, and the phosphate is released as inorganic phosphate [60].



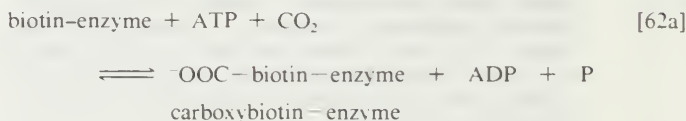
Lipid components. The component building blocks of the lipids found in storage fats, in lipoproteins (combinations of lipid and protein), and in the membranes of cells and organelles are glycerol, the fatty acids, and a number of other compounds (*e.g.*, serine, inositol).

Building blocks of lipids

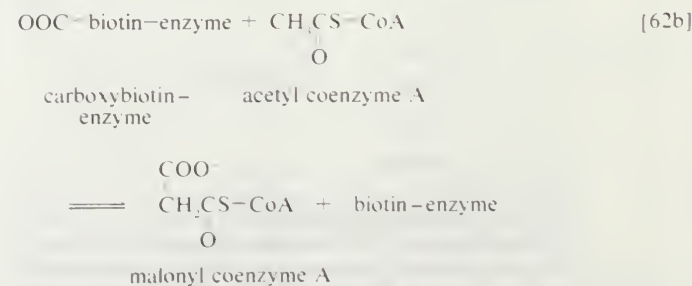
Glycerol. Glycerol is readily derived from dihydroxyacetone phosphate, an intermediate of glycolysis (see [4]). In a reaction catalyzed by glycerol 1-phosphate dehydrogenase [61], dihydroxyacetone phosphate is reduced to glycerol 1-phosphate; reduced NAD^+ provides the reducing equivalents for the reaction and is oxidized. This compound reacts further (see below *Other components*).



Fatty acids. Although all the carbon atoms of the fatty acids found in lipids are derived from the acetyl coenzyme A produced by the catabolism of carbohydrates and fatty acids (Figure 2), the molecule first undergoes a carboxylation, forming malonyl coenzyme A, before participating in fatty acid synthesis. The carboxylation reaction is catalyzed by acetyl CoA carboxylase, an enzyme whose prosthetic group is the vitamin biotin. The biotin-enzyme first undergoes a reaction that results in the attachment of carbon dioxide to biotin; ATP is required and forms ADP and inorganic phosphate [62a]. The complex product,

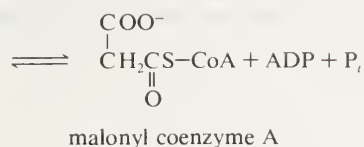
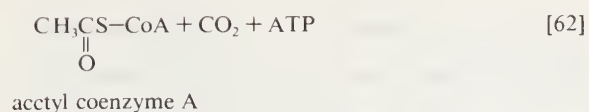


called carboxybiotin-enzyme, releases the carboxy moiety to acetyl coenzyme A, forming malonyl coenzyme A and restoring the biotin-enzyme [62b].



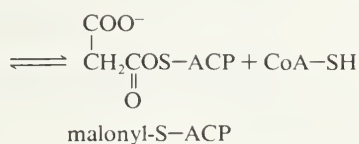
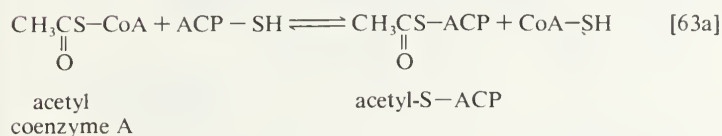
The overall reaction [62] catalyzed by acetyl coenzyme A carboxylase thus involves the expenditure of one molecule of ATP for the formation of each molecule of malonyl coenzyme A from acetyl coenzyme A and carbon dioxide.

Malonyl coenzyme A and a molecule of acetyl coenzyme A react (in bacteria) with the sulfhydryl group of

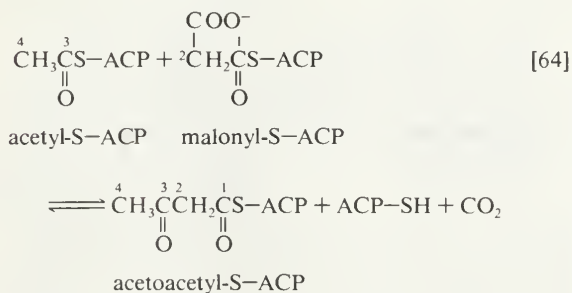


Acyl-carrier protein

a relatively small molecule known as acyl-carrier protein (ACP-SH); in higher organisms ACP-SH is part of a multi-enzyme complex called fatty acid synthetase. ACP-SH is involved in all of the reactions leading to the synthesis of a fatty acid such as palmitic acid from acetyl coenzyme A and malonyl coenzyme A. The products of [63a] and [63b] are acetyl-S-ACP, malonyl-S-ACP, and coenzyme A. The enzymes catalyzing [63a] and [63b] are known as acetyl transacylase and malonyl transacylase, respectively. Acetyl-ACP and malonyl-ACP react in a reaction catalyzed by β -ketoacyl-ACP synthetase so that the acetyl moiety ($\text{CH}_3\text{CO}-$) is transferred to the malonyl moiety ($-\text{OOCCH}_2\text{CO}-$). Simultaneously, the carbon dioxide fixed in step [62] is lost, leaving as a product a four-carbon moiety attached to ACP and called acetoacetyl-S-ACP [64].

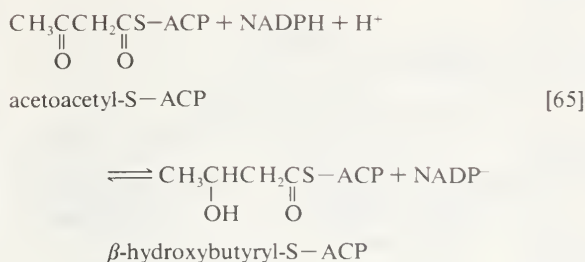


It should be noted that the carbon atoms of acetyl-S-ACP occur at the end of acetoacetyl-S-ACP (see carbon atoms numbered 4 and 3 in [64]) and that carbon dioxide plays an essentially catalytic role; the decarboxylation of the malonyl-S-ACP [64] provides a strong thermodynamic pull toward fatty acid synthesis.

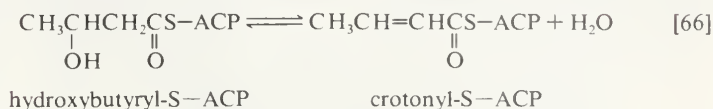


Fatty acid synthesis and breakdown

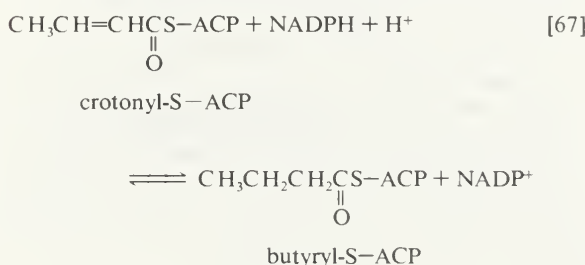
The analogy between reaction [64] of fatty acid synthesis and the cleavage reaction [25] of fatty acid catabolism is apparent in the other reactions of fatty acid synthesis. The acetoacetyl-S-ACP, for example, undergoes reduction to β -hydroxybutyryl-S-ACP [65]; the reaction is catalyzed by β -ketoacyl-ACP reductase. Reduced NADP⁺ is the electron donor, however, and not reduced NAD⁺ (which would participate in the reversal of reaction [24]). NADP⁺ is thus a product in [65]. In [66] β -hydroxybutyryl-S-ACP is dehydrated (*i.e.*, one molecule of water is removed), in a reaction catalyzed by enoyl-ACP-hydase, and then undergoes a second reduction [67], in which reduced NADP⁺



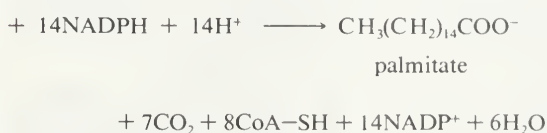
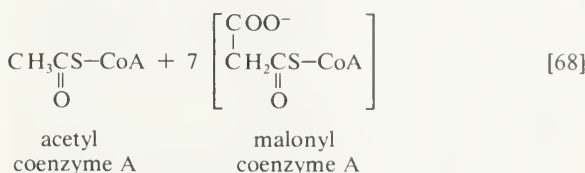
again acts as the electron donor. The products of [66] are crotonyl-S-ACP and water. The products of [67], which is catalyzed by crotonyl-ACP reductase, are butyryl-S-ACP and NADP⁺.



The formation of butyryl-S-ACP [67] completes the first of several cycles, in each of which one molecule of malonyl coenzyme A enters via reactions [62] and [63b]. In the cycle following the one ending with [67], the butyryl moiety is transferred to malonyl-S-ACP, and a molecule of carbon dioxide is again lost; a six-carbon compound results. In subsequent cycles, each of which adds two carbon atoms to the molecule via reaction [64], successively longer β -oxoacyl-S-ACP derivatives are produced.



Ultimately, a molecule with 16 carbon atoms, palmityl-S-ACP, is formed. In most organisms a deacylase catalyzes the release of free palmitic acid; in a few, synthesis continues, and an acid with 18 carbon atoms is formed. The fatty acids can then react with coenzyme A (compare reaction [21]) to form fatty acyl coenzyme A, which can condense with the glycerol 1-phosphate formed in step [61]; the product is a phosphatidic acid. The overall formation of each molecule of palmitic acid from acetyl coenzyme A—via step [62] and repeated cycles of steps [63] through [67]—requires the investment of seven molecules of ATP and 14 of reduced NADP⁺ (see [68]). The process is thus an energy-requiring one (endergonic) and represents a major way by which the reducing power generated in NADP-linked dehydrogenation reactions of carbohydrate catabolism is utilized (see above *The phosphogluconate pathway*).



Other components. The major lipids that serve as components of membranes, called phospholipids, as well as lipoproteins, contain, in addition to two molecules of fatty acid, one molecule of a variety of different compounds. The precursors of these compounds include serine, inositol, and glycerol 1-phosphate. They are derived from

Phospho-lipids

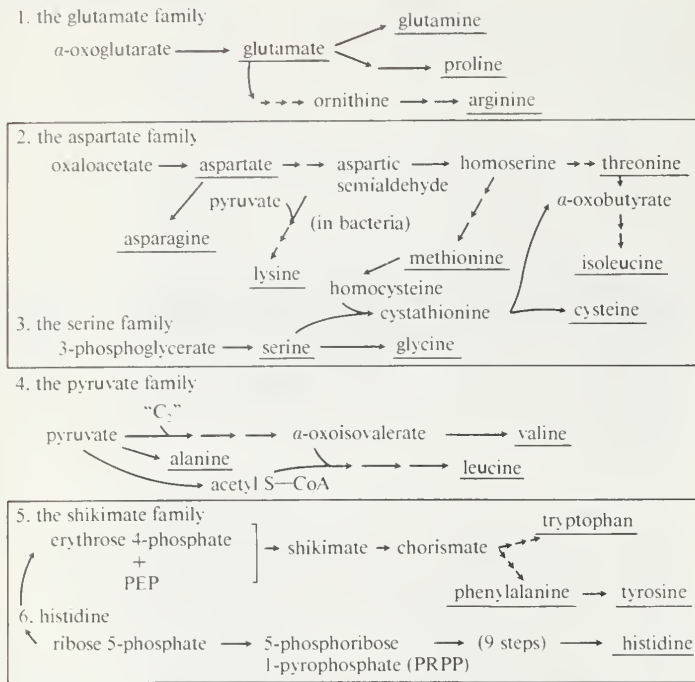


Figure 10: Family relationships in amino-acid biosyntheses. Components of proteins are underlined. Not all of the intermediates formed are named.

intermediates of the central metabolic pathways (e.g., Figure 10; reaction [61]).

Amino acids. Organisms differ considerably in their ability to synthesize amino acids from the intermediates of central metabolic pathways. Most vertebrates can form only the chemically most simple amino acids; the others must be supplied in the diet. Man, for example, synthesizes about 10 of the 20 commonly encountered amino acids; these are termed nonessential amino acids. The essential amino acids must be supplied in food.

Higher plants are more versatile than animals; they can make all of the amino acids required for protein synthesis, with either ammonia (NH₃) or nitrate (NO₃⁻) as the nitrogen source. Some bacteria, and leguminous plants (e.g., peas) that harbour such bacteria in their root nodules, are able to utilize nitrogen from the air to form ammonia and use the latter for amino-acid synthesis.

Bacteria differ widely in their ability to synthesize amino acids. Some species, such as *Escherichia coli*, which can grow in media supplied with only a single carbon source and ammonium salts, can make all of their amino acids from these starting materials. Other bacteria may require as many as 16 different amino acids.

Each of the 20 common amino acids is synthesized by a different pathway, the complexity of which reflects the chemical complexity of the amino acid formed. As with other compounds, the pathway for the synthesis of an amino acid is for the most part different from that by which it is catabolized. A detailed discussion of the pathway by which each amino acid is formed is beyond the scope of this article, but two salient features of amino-acid biosynthesis should be mentioned.

First, ammonia is incorporated into the intermediates of metabolic pathways mainly via the glutamate dehydrogenase reaction [28], which proceeds from right to left in biosynthetic reactions. Similarly, the transaminase enzymes (reactions [26a, b, and c]) enable the amino group (NH₂-) to be transferred to other amino acids.

Second, a group of several amino acids may be synthesized from one amino acid, which acts as a "parent" of an amino-acid "family." The families are also interrelated in several instances. Figure 10 shows, for bacteria that can synthesize 20 amino acids, the way in which they are derived from intermediates of pathways already considered. Alpha-oxoglutarate and oxaloacetate are intermediates of the TCA cycle; pyruvate, 3-phosphoglycerate, and PEP are intermediates of glycolysis; and ribose 5-phosphate

and erythrose 4-phosphate are formed in the phosphoglucuronate pathway.

Mononucleotides. Most organisms can synthesize the purine and pyrimidine nucleotides that serve as the building blocks of RNA (containing nucleotides in which the pentose sugar is ribose, called ribonucleotides) and DNA (containing nucleotides in which the pentose sugar is deoxyribose, called deoxyribonucleotides) as well as the agents of energy exchange.

Purine ribonucleotides. The purine ribonucleotides (AMP and GMP) are derived from ribose 5-phosphate. The overall sequence that leads to the parent purine ribonucleotide, which is inosinic acid, involves 10 enzymatic steps.

Figure 11 is an outline of the manner in which inosinic acid is synthesized. Inosinic acid can be converted to AMP and GMP; these in turn yield the triphosphates (i.e., ATP and GTP) via reactions catalyzed by adenylate kinase [69] and nucleoside diphosphate kinase (see reaction [43a]).

Inosinic acid



Pyrimidine ribonucleotides. The biosynthetic pathway for the pyrimidine nucleotides is somewhat simpler than that for the purine nucleotides. Aspartate (derived from

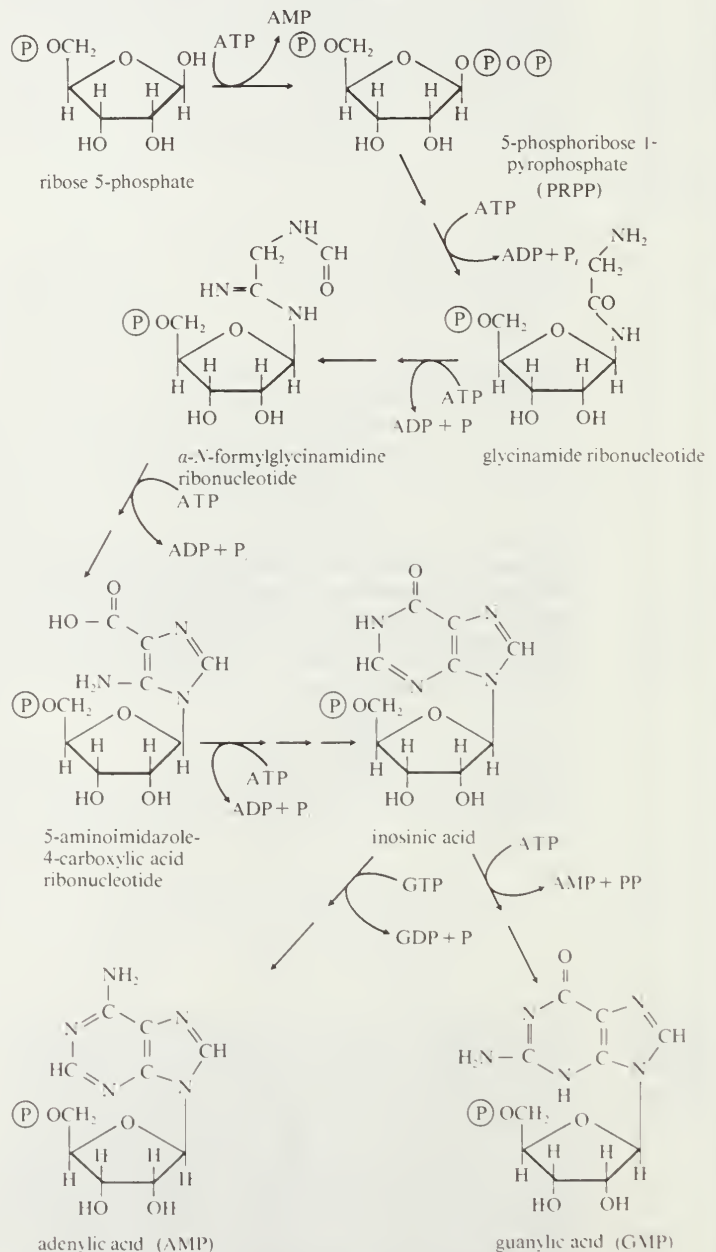
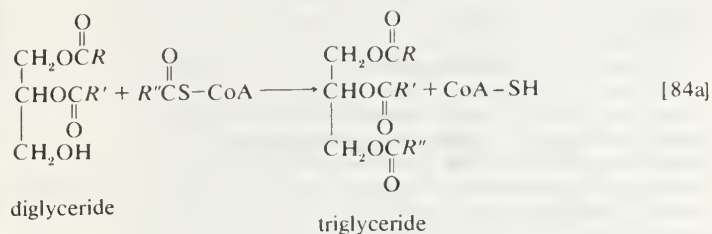
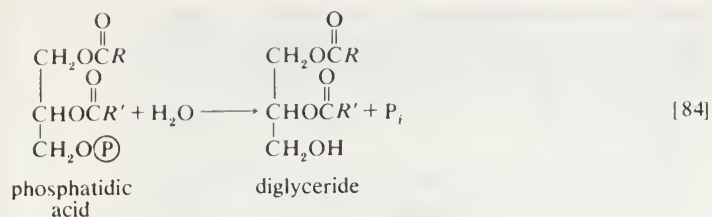


Figure 11: Biosynthesis of purine nucleotides. Not all of the intermediate compounds formed are shown.

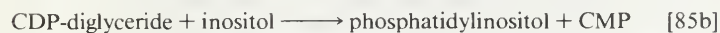
Amino-acid synthesis



phosphate group [84]; the diglyceride thus formed can then accept a third molecule of fatty acyl coenzyme A (represented as $\text{R}''\text{CS}-\text{CoA}$ in [84a]).

Bio-synthesis of phospholipids

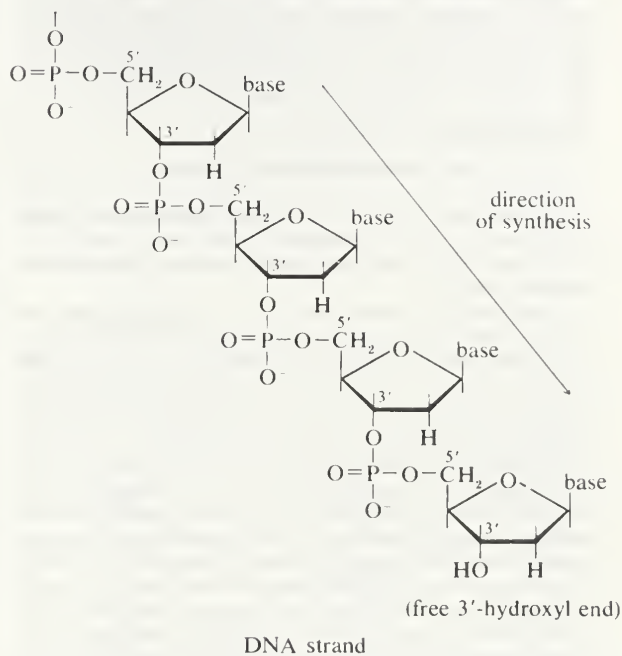
In the biosynthesis of phospholipids, however, phosphatidic acid is not hydrolyzed; rather, it acts as the $\text{R}-\oplus$ in reaction [77], the NTP here being cytidine triphosphate (CTP). A CDP-diglyceride is produced, and inorganic pyrophosphate is released [77b]. CDP-diglyceride is the common precursor of a variety of phospholipids. In subsequent reactions, each catalyzed by a specific enzyme, CMP is displaced from CDP-diglyceride by one of three compounds—serine, inositol, or glycerol 1-phosphate—to form CMP and, respectively, phosphatidylserine [85a], phosphatidylinositol [85b], or, in [85c], 3-phosphatidyl-glycerol 1'-phosphate (PGP). These reactions differ from those of polysaccharide biosynthesis ([79], [82]) in that phosphate is retained in the phospholipid, and the nucleotide product (CMP) is therefore a nucleoside monophosphate rather than the diphosphate. These compounds can react further: phosphatidylserine to give, sequentially, phosphatidylethanolamine and phosphatidylcholine; phosphatidylinositol to yield mono- and diphosphate derivatives that are components of brain tissue and of mitochondrial membranes; and PGP to yield the phosphatidylglycerol abundant in many bacterial membranes and the diphosphatidylglycerol that is also a major component of mitochondrial and bacterial membranes.



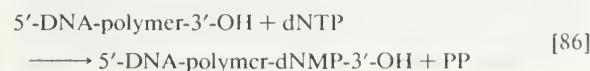
Formation of nucleic acids

Nucleic acids and proteins. As with the synthesis of polysaccharides and lipids, the formation of the nucleic acids and proteins from their building blocks requires the input of energy. Nucleic acids are formed from nucleoside triphosphates, with concomitant elimination of inorganic pyrophosphate, which is subsequently hydrolyzed via reaction [21a]. Amino acids also are activated, forming, at the expense of ATP, aminoacyl-complexes. This activation process is also accompanied by loss of inorganic pyrophosphate. But, although these biochemical processes are basically similar to those involved in the biosynthesis of other macromolecules, their occurrence is specifically subservient to the genetic information in DNA. DNA contains within its structure the blueprint both for its own exact duplication and for the synthesis of a number of types of RNA, among which is a class termed messenger RNA (mRNA). A complementary relationship exists between the sequence of purines (*i.e.*, adenine and guanine) and pyrimidines (cytosine and thymine) in the DNA comprising a gene and the sequence in mRNA into which this genetic information is transcribed. This information is then translated into the sequence of amino acids in a protein, a process that involves the functioning of a variety of other classes of ribonucleic acids (see GENETICS AND HEREDITY, THE PRINCIPLES OF).

Synthesis of DNA. The maintenance of genetic integrity demands not only that enzymes exist for the synthesis of DNA but that they function so as to ensure the replication of the genetic information (encoded in the DNA to be copied) with absolute fidelity. This implies that the assembly of new regions of a DNA molecule must occur on a template of DNA already present in the cell. The synthetic processes must also be capable of repairing limited regions of DNA, which may have been damaged, for example, as a consequence of exposure to ultraviolet irradiation. The physical structure of DNA is ideally adapted to its biological roles. Two strands of nucleotides are wound around each other in the form of a double helix. The helix is stabilized by hydrogen bonds that occur between the purine and pyrimidine bases of the strands. Thus, the adenine of one strand pairs with the thymine of the other, and the guanine of one strand with the cytosine of the other. The base pairs may be visualized as the treads of a spiral staircase, in which the two chains of repeating units (*i.e.*, ribose-phosphate-ribose) form the sides.



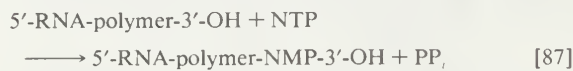
During the biosynthesis of DNA, the two strands unwind, and each serves as a template for the synthesis of a new, complementary strand, in which the bases pair in exactly the same manner as occurred in the parent double helix. The process is catalyzed by a DNA polymerase enzyme, which catalyzes the addition of the appropriate deoxyribonucleoside triphosphate (NTP) in [86] onto one end, specifically, the free 3'-hydroxyl end ($-\text{OH}$) of the growing DNA chain (see diagram of DNA strand). In [86] the addition of a deoxyribonucleoside monophosphate (dNMP) moiety onto a growing DNA chain ($5'$ -DNA-polymer- $3'$ -OH) is shown; the other product is inorganic pyrophosphate. The specific nucleotide inserted in the growing chain is dictated by the base in the complementary (template) strand of DNA with which it pairs. The functioning of DNA polymerase thus requires the presence of all four deoxyribonucleoside triphosphates (*i.e.*, dATP, dTTP, dGTP, and dCTP) as well as preformed DNA to act as a template. Although a number of DNA polymerase enzymes have been purified from different organisms, it is not yet certain whether those that have been most extensively studied are necessarily involved in the formation of new DNA molecules, or whether they are primarily concerned with the repair of damaged regions of molecules. A polynucleotide ligase that effects the formation of the phosphate bond between adjacent sugar molecules is concerned with the repair function but may also have a role in synthesis.



Bio-synthesis of DNA

The varieties of RNA

Synthesis of RNA. Various types of RNA are found in living organisms: messenger RNA (mRNA) is involved in the immediate transcription of regions of DNA; transfer RNA (tRNA) is concerned with the incorporation of amino acids into proteins; and structural RNA is found in the ribosomes that form the protein-synthesizing machinery of the cell. In cells of organisms with well-defined nuclei (*i.e.*, eukaryotes), a heterogeneous RNA fraction of unknown function is constantly broken down and resynthesized in the nucleus of the cell but does not leave it. The different types of RNA are synthesized via RNA polymerases [87], the action of which is analogous to that of the DNA polymerases that catalyze reaction [86]. In [87] the growing RNA chain is represented by 5'-RNA-polymer-3'-OH, and the ribonucleoside triphosphate by NTP. One product (5'-RNA-polymer-NMP-3'-OH) reflects the incorporation of ribonucleoside monophosphate; the other product is, as in [86], inorganic pyrophosphate. Synthesis of RNA requires DNA as a template, thus ensuring that the base composition of the RNA faithfully reflects that of the DNA; in addition, as in DNA synthesis, all four nucleoside triphosphates must be present. The major differences between reactions [86] and [87] are that, in the latter, the nucleotides contain ribose instead of deoxyribose, and that, in RNA, uracil replaces the thymine of DNA.



It appears that, although only one strand of the DNA double helix serves as template during the formation of RNA, some regions are transcribed from one strand, some from the other.

An important constraint on RNA synthesis is that the accurate copying of the appropriate DNA strand by RNA polymerase must start at the beginning of a gene—and not somewhere along it—and must stop as soon as the genetic information has been transcribed. The way in which this selectivity is achieved is not yet fully understood, although it has been established that *E. coli* contains a protein, the sigma factor, that is not required for the incorporation of the nucleoside triphosphates into the growing RNA chain but apparently is essential for binding RNA polymerase to the proper DNA sites to initiate RNA synthesis. After the initiation step, the sigma factor is released; the role of the sigma factor in transcription suggests that the DNA at the initiation sites must be unique in some way so as to ensure that the correct strand is used as the template. Evidence indicates further that other protein factors are involved in the termination of transcription.

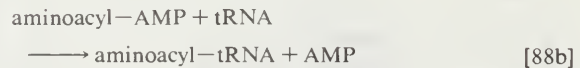
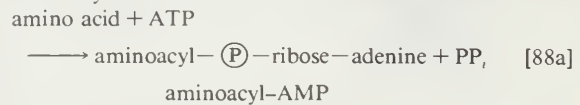
Synthesis of proteins. Approximately 120 macromolecules are involved directly or indirectly in the process of the translation of the base sequence of a messenger RNA molecule into the amino-acid sequence of a protein. The relationship between the base sequence and the amino-acid sequence constitutes the genetic code. The basic properties of the code are: it is triplet—*i.e.*, a linear sequence of three bases in mRNA specifies one amino acid in a protein; it is nonoverlapping—*i.e.*, each triplet is discrete and does not overlap either neighbour; it is degenerate—*i.e.*, many of the 20 amino acids are specified by more than one of the 64 possible triplets of bases; and it appears to apply universally to all living organisms.

The main sequence of events associated with the expression of this genetic code, as elucidated for *E. coli*, may be summarized as follows (see also GENETICS AND HEREDITY, THE PRINCIPLES OF).

1. Messenger RNA binds to the smaller of two subunits of large particles termed ribosomes.

2. The amino acid that begins the assembly of the protein chain is activated and transferred to a specific transfer RNA (tRNA). The activation step, catalyzed by an aminoacyl-tRNA synthetase specific for a particular amino acid, effects the formation of an aminoacyl-AMP complex [88a] in a manner somewhat analogous to reaction [77]; ATP is required, and inorganic pyrophosphate is a product. The aminoacyl-AMP, which remains bound to the enzyme, is transferred to a specific molecule of tRNA in a reaction catalyzed by the same enzyme. AMP is re-

leased, and the other product is called aminoacyl-tRNA [88b]. In *E. coli* the amino acid that begins the assembly of the protein is always formylmethionine (f-Met). There is no evidence that f-Met is involved in protein synthesis in eukaryotic cells.



3. Aminoacyl-tRNA binds to the mRNA-ribosomal complex in a reaction in which energy is provided by the hydrolysis of GTP to GDP and inorganic phosphate. In this step and in 5 below, the genetic code is translated. All of the different tRNAs contain triplets of bases that pair specifically with the complementary base triplets in mRNA; the base triplets in mRNA specify the amino acids to be added to the protein chain. During or shortly after the pairing occurs the aminoacyl-tRNA moves from the aminoacyl-acceptor (A) site on the ribosome to another site, called a peptidyl-donor (P) site.

4. The larger subunit of the ribosome then joins the mRNA-f-Met-tRNA-smaller ribosomal subunit complex.

5. The second amino acid to be added to the protein chain is specified by the triplet of bases adjacent to the initiator triplet in mRNA. The amino acid is activated and transferred to its tRNA by a repetition of reactions [88a] and [88b]. This newly formed aminoacyl-tRNA now binds to the A site of the mRNA-ribosome complex, with concomitant hydrolysis of GTP.

6. The enzyme peptidyl transferase, which is part of the larger of the two ribosomal subunits, catalyzes the transfer of formylmethionine from the tRNA to which it is attached (designated tRNA^{f-Met}) to the second amino acid; for example, if the second amino acid were leucine, step 5 would have achieved the binding of leucyl-tRNA (Leu-tRNA^{Leu}) next to f-Met-tRNA^{f-Met} on the ribosome-mRNA complex. Step 6 catalyzes the transfer reaction that is shown in [89], in which tRNA^{f-Met} is released from formyl-methionine (f-Met), and Leu-tRNA^{Leu} is bound to formyl-methionine.



7. In the next step three results are achieved. The dipeptide f-Met-Leu (a dipeptide consists of two amino acids) moves from the A (aminoacyl-acceptor) site to the P (peptidyl-donor) site on the ribosome; the tRNA^{f-Met} is thereby displaced from the P site, and the ribosome moves the length of one triplet (three bases) along the mRNA molecule. The occurrence of these events is accompanied by the hydrolysis of a second molecule of GTP and leaves the system ready to receive the next aminoacyl-tRNA (by repetition of step 5). The cycle of events in 5, 6, and 7 is repeated until the ribosome moves to a triplet on the mRNA that does not specify an amino acid but provides the signal for termination of the amino-acid chain. Triplets of this type are represented by one uracil (U) preceding, and adjacent to, two adenines (UAA) or preceding one adenine and one guanosine in either order (UGA, or UAG).

8. At the termination of synthesis the completed protein is released from the tRNA to which it had remained linked. Two further events then occur in *E. coli*. First, the formyl constituent of the f-methionyl moiety is hydrolyzed by the catalytic action of a formylase, producing a protein with methionine at the end. If the required protein does not contain methionine in this position (and the majority of proteins in *E. coli* appear to), the methionine and possibly other amino acids that follow it are removed by enzymatic reactions. Second, the ribosome-mRNA complex dissociates, and the ribosomal subunits become available for a new round of translation by binding another mRNA molecule, step 1.

Translation

Termination of the amino-acid chain

For the sake of brevity, other ancillary protein factors that participate in this sequence 1 to 8 have been omitted; the role of many of these factors is as yet poorly understood.

Regulation of metabolism

FINE CONTROL

The flux of nutrients along each metabolic pathway is governed chiefly by two factors: (1) the availability of substrates on which pacemaker, or key, enzymes of the pathway can act and (2) the intracellular levels of specific metabolites that affect the reaction rates of pacemaker enzymes. Key enzymes are usually complex proteins that, in addition to the site at which the catalytic process occurs (*i.e.*, the active site), contain sites to which the regulatory metabolites bind. Interactions with the appropriate molecules at these regulatory sites cause changes in the shape of the enzyme molecule. Such changes may either facilitate or hinder the changes that occur at the active site. The rate of the enzymatic reaction is thus speeded up or slowed down by the presence of a regulatory metabolite.

In many cases, the specific small molecules that bind to the regulatory sites have no obvious structural similarity to the substrates of the enzymes; these small molecules are therefore termed allosteric effectors, and the regulatory sites are termed allosteric sites (see the article **BIOCHEMICAL COMPONENTS OF ORGANISMS**). Allosteric effectors may be formed by enzyme-catalyzed reactions in the same pathway in which the enzyme regulated by the effectors functions. In this case a rise in the level of the allosteric effector would affect the flux of nutrients along that pathway in a manner analogous to the feedback phenomena of homeostatic processes. Such effectors may also be formed by enzymatic reactions in apparently unrelated pathways. In this instance the rate at which one metabolic pathway operates would be profoundly affected by the rate of nutrient flux along another. It is this situation that, to a large extent, governs the sensitive and immediately responsive coordination of the many metabolic routes in the cell.

End-product inhibition. A biosynthetic pathway is usually controlled by an allosteric effector produced as the end product of that pathway, and the pacemaker enzyme on which the effector acts usually catalyzes the first step that uniquely leads to the end product. This phenomenon, called end-product inhibition, is illustrated by the multienzyme, branched pathway for the formation from oxaloacetate of the "aspartate family" of amino acids (Figure 10). The system of interlocking controls is described in greater detail in Figure 12. As mentioned previously in this article, only plants and microorganisms can synthesize many of these amino acids, most animals requiring such amino acids to be supplied preformed in their diets.

Figure 12 shows that there are a number of pacemaker enzymes in the biosynthetic route for the aspartate family of amino acids, most of which are uniquely involved in the formation of one product. Each of the enzymes functions after a branch point in the pathway, and all are inhibited specifically by the end product that emerges from the branch point. It is not difficult to visualize from Figure 12 how the supplies of lysine, methionine, and isoleucine required by a cell can be independently regulated. Threonine, however, is both an amino acid essential for protein synthesis and a precursor of isoleucine. If the rate of synthesis of threonine from aspartate were regulated as are the rates of lysine, methionine, and isoleucine, an imbalance in the supply of isoleucine might result. This risk is overcome in *E. coli* by the existence of three different aspartokinase enzymes, all of which catalyze the first step common to the production of all the products derived from aspartate. Each has a different regulatory effector molecule. Thus, one type of aspartokinase is inhibited by lysine, a second by threonine. The third kinase is not inhibited by any naturally occurring amino acid, but its rate of synthesis (see below) is controlled by the concentration of methionine within the cell. The triple control mechanism resulting from the three different aspartokinases ensures that the accumulation of one amino acid does not shut off the supply of aspartyl phosphate necessary for the synthesis of the others.

Another example of control through end-product inhibition also illustrates the manner in which the operation of two biosynthetic pathways may be coordinated. Both DNA and the various types of RNA are assembled from purine and pyrimidine nucleotides (see above *Nucleic acids and proteins*); these in turn are built up from intermediates of central metabolic pathways (see above *Mononucleotides*). The first step in the synthesis of pyrimidine nucleotides is that catalyzed by aspartate carbamoyltransferase [70a]. This step initiates a sequence of reactions that leads to the formation of pyrimidine nucleotides such as UTP and CTP [74]. Studies of aspartate carbamoyltransferase have revealed that the affinity of this enzyme for its substrate (aspartate) is markedly decreased by the presence of CTP. This effect can be overcome by the addition of ATP, a purine nucleotide. The enzyme can be dissociated into two subunits: one contains the enzymatic activity and (in the dissociated form) does not bind CTP; the other binds CTP but has no catalytic activity. Apart from providing physical evidence that pacemaker enzymes contain distinct catalytic and regulatory sites, the interaction of aspartate carbamoyltransferase with the different nucleotides provides an explanation for the control of the supply of nucleic acid precursors. If a cell contains sufficient pyrimidine nucleotides (*e.g.*, UTP), aspartate carbamoyltransferase, the first enzyme of pyrimidine biosynthesis, is inhibited. If, however, the cell contains high levels of purine nucleotides (*e.g.*, ATP), as required for the formation of nucleic acids, the inhibition of aspartate carbamoyltransferase is relieved, and pyrimidines are formed.

Aspartate carbamoyltransferase

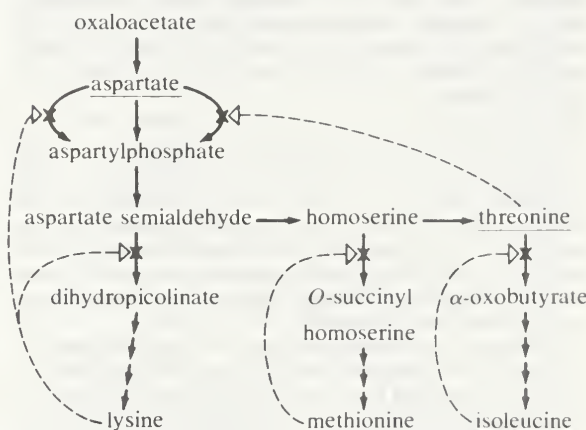


Figure 12: "Fine control" of the enzymes of the aspartate family in *E. coli* (see text).

Positive modulation. Not all pacemaker enzymes are controlled by inhibition of their activity. Instead, some are subject to positive modulation—*i.e.*, the effector is required for the efficient functioning of the enzyme. Such enzymes exhibit little activity in the absence of the appropriate allosteric effector. One instance of positive modulation is the anaplerotic fixation of carbon dioxide onto pyruvate and phosphoenolpyruvate (PEP); this example also illustrates how a metabolic product of one route controls the rate of nutrient flow of another (see Figure 9).

The carboxylation of pyruvate in higher organisms [50] and the carboxylation of phosphoenolpyruvate in gut bacteria [50a] occurs at a significant rate only if acetyl coenzyme A is present. Acetyl coenzyme A acts as a positive allosteric effector and is not broken down in the course of the reaction. Moreover, some pyruvate carboxylases [50] and the PEP carboxylase of gut bacteria are inhibited by four-carbon compounds (*e.g.*, aspartate). These substances inhibit because they interfere with the binding of the positive effector, acetyl coenzyme A. Such enzymatic controls are reasonable in a physiological sense: it will be recalled that anaplerotic formation of oxaloacetate from pyruvate or PEP is required to provide the acceptor for the entry of acetyl coenzyme A into the TCA cycle. The reaction need occur only if acetyl coenzyme A is present in sufficient amounts. On the other hand, an abundance of four-carbon intermediates obviates the necessity for forming more through carboxylation reactions such as [50] and [50a].

Rate of the
glyoxylate
cycle

Similar reasoning, though in the opposite sense, can be applied to the control of another anaplerotic sequence, the glyoxylate cycle (Figure 8). The biosynthesis of cell materials from the two-carbon compound acetate is, in principle, akin to biosynthesis from TCA cycle intermediates. In both processes, it is the availability of intermediates such as PEP and pyruvate that determines the rate at which a cell forms the many components produced through gluconeogenesis. Although in the strictest sense the glyoxylate cycle has no defined end product, PEP and pyruvate are, for these physiological reasons, best fitted to regulate the rate at which the glyoxylate cycle is required to operate. It is thus not unexpected that the pacemaker enzyme of the glyoxylate cycle, isocitrate lyase (reaction [52]), is allosterically inhibited by PEP and by pyruvate.

Energy state of the cell. It is characteristic of catabolic routes that they do not lead to uniquely identifiable end products. The major products of glycolysis and the TCA cycle, for example, are carbon dioxide and water. Within the cell, the concentrations of both are unlikely to vary sufficiently to allow them to serve as effective regulatory metabolites. The processes by which water is produced (Figure 7) initially involve, however, the reduction of coenzymes, the reoxidation of which is accompanied by the synthesis of ATP from ADP. Moreover, as described in previous sections, the utilization of ATP in energy-consuming reactions yields ADP and AMP. At any given moment, therefore, a living cell contains ATP, ADP, and AMP; the relative proportion of the three nucleotides provides an index of the energy state of the cell. It is thus reasonable that the flux of nutrients through catabolic routes is, in general, impeded by high intracellular levels of both reduced coenzymes (*e.g.*, FADH₂, reduced NAD⁺) and ATP, and that these inhibitory effects are often overcome by AMP.

The control exerted by the levels of ATP, ADP, and AMP within the cell is illustrated by the regulatory mechanisms of glycolysis and the TCA cycle (Figure 9); these nucleotides also serve to govern the occurrence of the opposite pathway, gluconeogenesis, and to avoid mutual interference of the catabolic and anabolic sequences. Although not all of the controls mentioned below have been found to operate in all living organisms examined, it has been observed that, in general:

1. Glucose 6-phosphate stimulates glycogen synthesis from glucose 1-phosphate [79] and inhibits both glycogen breakdown [16] and its own formation from glucose [1].

2. Phosphofructokinase, the most important pacemaker enzyme of glycolysis [3], is inhibited by high levels of its own substrates (fructose 6-phosphate and ATP); this inhibition is overcome by AMP. In tissues, such as heart muscle, which use fatty acids as a major fuel, inhibition of glycolysis by citrate may be physiologically the more important means of control. Control by citrate, the first intermediate of the TCA cycle, which produces the bulk of the cellular ATP, is thus the same, in principle, as control through ATP.

3. Fructose 1,6-diphosphatase [59], which catalyzes the reaction opposite to phosphofructokinase, is strongly inhibited by AMP.

4. Rapid catabolism of carbohydrate requires the efficient conversion of PEP to pyruvate. In liver and in some bacteria the activity of the pyruvate kinase that catalyzes this process [10] is greatly stimulated by the presence of fructose 1,6-diphosphate, which thus acts as a potentiator of a reaction required for its ultimate catabolism.

Inhibition
of pyruvate
oxidation

5. The oxidation of pyruvate to acetyl coenzyme A [37] is inhibited by acetyl coenzyme A. Because acetyl coenzyme A also acts as a positive modulator of pyruvate carboxylation [50], this control reinforces the partition between pyruvate catabolism and its conversion to four-carbon intermediates for anaplerosis and gluconeogenesis.

6. Citrate synthase [38], the first enzyme of the TCA cycle, is inhibited by ATP in higher organisms and by reduced NAD⁺ in many microorganisms. In some strictly aerobic bacteria, the inhibition by reduced NAD⁺ is overcome by AMP.

7. Citrate acts as a positive effector for the first enzyme of fatty acid biosynthesis [62]. A high level of citrate,

which also indicates a sufficient energy supply, thus inhibits carbohydrate fragmentation (see [2]) and diverts the carbohydrate that has been fragmented from combustion to the formation of lipids.

8. Some forms of isocitrate dehydrogenase [40] are maximally active only in the presence of ADP or AMP and are inhibited by ATP. This is an example of regulation by covalent modification of an enzyme since the action of ATP here is to phosphorylate, and consequently to inactivate, the isocitrate dehydrogenase. A specific phosphatase, which is a different enzymatic activity of the protein that effects the phosphorylation by ATP, catalyzes the splitting-off by water of the phosphate moiety on the inactive isocitrate dehydrogenase and thus restricts activity. Again, the energy state of the cell serves as the signal regulating an enzyme involved in energy transduction.

COARSE CONTROL

Although fine control mechanisms allow the sensitive adjustment of the flux of nutrients along metabolic pathways relative to the needs of cells under relatively constant environmental conditions, these processes may not be adequate to cope with severe changes in the chemical milieu.

Such severe changes may arise in higher organisms with a change in diet or when, in response to other stimuli, the hormonal balance is altered (see **BIOCHEMICAL COMPONENTS OF ORGANISMS**). In starvation, for example, the overriding need to maintain blood glucose levels may require the liver to synthesize glucose from noncarbohydrate products of tissue breakdown at rates greater than can be achieved by the enzymes normally present in the liver. Under such circumstances, cellular concentrations of key enzymes of gluconeogenesis, such as pyruvate carboxylase [50] and PEP carboxykinase [54], may rise by as much as 10-fold, while the concentration of glucokinase [1] and of the enzymes of fatty acid synthesis decreases to a similar extent. Conversely, high carbohydrate diets and administration of the hormone insulin to diabetic animals elicit a preferential synthesis of glucokinase [1] and pyruvate kinase [10]. These changes in the relative proportions and absolute amounts of key enzymes are the net result of increases in the rate of their synthesis and decreases in the rate of their destruction. Although such changes reflect changes in the rates of either transcription, translation, or both of specific regions of the genome, the mechanisms by which the changes are effected have not yet been clarified.

Microorganisms sometimes encounter changes in environment much more severe than those encountered by the cells of tissues and organs, and their responses are correspondingly greater. Mention has already been made of the ability of *E. coli* to form β -galactosidase when transferred to a medium containing lactose as the sole carbon source (see *Integration of catabolism and anabolism*): such a transfer may result in an increase of 1,000-fold or more in the cellular concentration of the enzyme. Because this preferential enzyme synthesis is elicited by exposure of the cells to lactose, or to non-metabolizable but chemically similar analogues, and because synthesis ceases as soon as the eliciting agents (inducers) are removed, β -galactosidase is termed an inducible enzyme. It has been established that a regulator gene exists that specifies the amino-acid sequence of a so-called repressor protein, and that the repressor protein binds to a unique portion of the region of DNA concerned with β -galactosidase formation. Under these circumstances the DNA is not transcribed to mRNA, and virtually no enzyme is made. The repressor, however, is an allosteric protein and readily combines with inducers. Such a combination prevents the repressor from binding to DNA and allows transcription and translation of β -galactosidase to proceed.

β -galactosi-
dase

Although this mechanism for the specific control of gene activity may not apply to the regulation of all inducible enzymes—for example, those concerned with the utilization of the sugar arabinose—and is not universally applicable to all coarse control processes in all microorganisms, it can explain the manner in which the presence in growth media of at least some cell components represses (*i.e.*, inhibits the synthesis of) enzymes normally involved in the formation of such components by gut bacteria such

as *E. coli*. Although, for example, the bacteria must obviously make amino acids from ammonia if that is the sole source of nitrogen available to them, it would not be necessary for the bacteria to synthesize enzymes required for the formation of amino acids supplied preformed in the medium. Thus, of the three aspartokinases formed by *E. coli* (Figure 12), two are repressed by their end products, methionine and lysine. On the other hand, the third aspartokinase, which (as described above) is inhibited by threonine, is repressed by threonine only if isoleucine is also present. This example of so-called multivalent repression is of obvious physiological utility. It is likely that the amino acids that thus specifically inhibit the synthesis of aspartokinases do so by combining with specific protein repressor molecules; however, whereas the combination of the inducer with the repressor of β -galactosidase inactivates the repressor protein and hence permits synthesis of the enzyme, the repressor proteins for biosynthetic enzymes would not bind to DNA unless they were also combined with the appropriate amino acid. Aspartokinase synthesis would thus occur in the absence of the end-product effectors and not in their presence.

This explanation applies also to the coarse control of the anaplerotic glyoxylate cycle (Figure 8). The synthesis of both of the enzymes unique to that cycle, isocitrate lyase [52] and malate synthase [53], is controlled by a regulator gene that presumably specifies a repressor protein unable to bind to DNA unless combined with pyruvate or PEP. Cells growing on acetate do not contain high levels of these intermediates because they are continuously being removed for biosynthesis. The enzymes of the glyoxylate cycle are therefore formed at high rates. If pyruvate or substances catabolized to PEP or pyruvate are added to the medium, however, further synthesis of the two enzymes is speedily repressed. (Ha.Ko.)

A regulator gene

Metabolic diseases

GENERAL CONSIDERATIONS

Metabolic diseases are conditions caused by an abnormality in one or more of the chemical reactions essential to producing energy, to regenerating cellular constituents, or to eliminating unneeded products arising from these processes. Depending on which metabolic pathway is involved, a single defective chemical reaction may produce consequences that are narrow, involving a single body function, or broad, affecting many organs and systems.

In addition to the fine and coarse mechanisms that affect individual chemical reactions (see above *Regulation of metabolism*), the rate and direction of metabolic processes are regulated by the availability of substrates (particularly essential amino acids or fatty acids), cofactors (whose presence is required to assist the enzyme), and key minerals; and by the influence of hormones secreted by endocrine glands and tissues. Abnormalities in any of these areas affect the many metabolic pathways that are driven by the catalytic activities of different enzymes. Nutritional and endocrine diseases are, in fact, the most common causes of metabolic dysfunction (see ENDOCRINE SYSTEMS), and diseases involving thyroid hormones and insulin, as well as caloric deprivation, illustrate the impact such abnormalities may have on the overall metabolic cycle.

The thyroid gland secretes two hormones, thyroxine and triiodothyronine, that regulate the overall rate of oxygen utilization and hence the production of useful energy from the combustion of food. When an excess of thyroid hormone is present (hyperthyroidism), the rate of oxidation of sugars, fatty acids, and amino acids is above normal. Simultaneously, as more oxygen is used, more carbon dioxide (CO₂) and urea are formed as the end products of metabolism. The demands for calories to supply the extra substrates and of vitamins as enzyme cofactors increase. Despite this, there is usually a loss of body weight from both fat and muscle tissue, the latter composed primarily of protein. Much of the extra energy generated is in the form of heat, and perspiration and evaporation of water are increased in order to maintain normal body temperature. Because the body requires more oxygen, the rates of respiration and blood flow increase to allow the ex-

tra oxygen to be delivered to tissues and to remove the extra carbon dioxide. A deficiency of thyroid hormone (hypothyroidism) produces exactly the opposite effects.

Another major hormone that influences metabolism is insulin, which is secreted by special cells of the pancreas. Insulin primarily regulates the direction of metabolism, shifting many processes toward the storage of substrates and away from their degradation. Insulin acts to increase the transport of glucose and amino acids as well as key minerals such as potassium, magnesium, and phosphate from the blood into cells. It also regulates a variety of enzymatic reactions within the cells, all of which have a common overall direction, namely the synthesis of large molecules from small units. A deficiency in the action of insulin (diabetes mellitus) causes severe impairment in (1) the storage of glucose in the form of glycogen and the oxidation of glucose for energy; (2) the synthesis and storage of fat from fatty acids and their precursors and the completion of fatty-acid oxidation; and (3) the synthesis of proteins from amino acids. As a result, high levels of glucose, fatty acids, keto acid products of incomplete fatty-acid oxidation, and essential amino acids are present in the blood. (The consequences of these abnormalities in metabolism are described in ENDOCRINE SYSTEMS.) Replacements for deficient hormones and methods to eliminate or combat excess hormones have been developed.

Insulin

A fundamental example of a nutritional cause for metabolic disease is caloric deprivation. In the most extreme case of fasting, the average human can survive about 60 days by oxidizing body stores of energy. Extremely obese persons can survive up to a year without food because of greater stores of fat and protein. During fasting, blood levels of glucose are low and the oxidation of glucose is curtailed, even by the brain tissue for which it is required. The blood levels of fatty acids and of their products of combustion, keto acids, are raised; keto acids partially take the place of glucose as a substrate in the brain. As the proteins of organs are degraded and not replaced, organ function declines, eventually resulting in death. Partial deprivation of calories results in a more gradual weight loss, a decline in oxygen use and energy production, a decrease in muscle strength, and lassitude. In infants and children, growth is impaired, and they are susceptible to infection and early death. In those children in whom calories are adequate but protein is chronically deficient (kwashiorkor), growth virtually ceases and the synthesis of important proteins, such as albumin, is reduced. Low levels of protein in the blood allow water to accumulate in tissues, causing the face, extremities, and abdomen to swell (edema). The immune system produces fewer cells to combat invading organisms, leading to infection and possibly death. Less severe problems include changes in skin and hair colour and texture and delayed sexual maturation. Treatment with protein or with the missing amino acids reverses the disorder.

THE DERIVATION OF SPECIFIC METABOLIC DISORDERS

Most metabolic diseases are caused by single defects in particular biochemical pathways, defects that are due to the deficient activity of individual enzymes. Around the turn of the century a British physician, Sir Archibald Garrod, proposed that such diseases were passed through generations, and he labeled them "inborn errors of metabolism." Subsequent studies have demonstrated the mechanisms by which a biochemical defect can be inherited. Beadle and Taum studied mutations in fruit flies and microorganisms and demonstrated that a single gene controls a single type of biochemical reaction. An American chemist, Linus Pauling, and a British biochemist, Vernon Ingram, showed that a mutation may be expressed as a specific alteration in the primary structure of a single type of protein.

Inborn errors of metabolism

Each enzyme is a protein and, like other proteins, has a unique sequence of amino acids. The order of the amino acids is directed by a segment of the DNA molecule called a gene. Within the gene, a sequence of three nucleotide bases directs each amino acid. If the nucleotide bases of the gene are rearranged, deleted, or substituted, the mutant gene directs the synthesis of a mutant enzyme protein, which may have only one amino acid substituted or large

segments of amino acids rearranged. Changes in structure usually decrease or entirely eliminate an enzyme's catalytic activity, slowing the rate of the biochemical reaction catalyzed by the enzyme. As a result, the amount of the product of the reaction is low, while that of the starting material (substrate) is high (Figure 13).

The potential consequences of mutant enzymes on cellular or organ function are varied: (1) the product of the affected metabolic pathway may have a key role that is not fulfilled; (2) a metabolic pathway that is repressed by normal levels of the product may proceed, with toxic consequences, when that product is at low levels; (3) excess substrate from the affected reaction may be toxic; (4) the excess substrate may become involved in and move rapidly through an alternate metabolic pathway, and the product of this pathway may be toxic. About 170 errors of metabolism can be traced to mutations in the genetic code, and the consequences range from trivial to lethal.

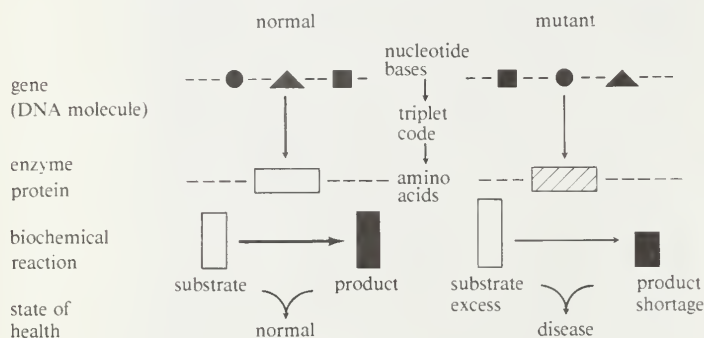


Figure 13: Pathological sequence of inborn errors of metabolism.

A change in the triplet nucleotide base sequence occurs at one point of the mutant gene, and the altered sequence codes for the "wrong" amino acid at a critical point in the structure of the enzyme protein. Because the mutant enzyme has little or no catalytic activity, its biochemical reaction is greatly diminished, causing a shortage of the product and an accumulation of the substrate.

The nucleus holds the DNA molecules in an amorphous strand called chromatin. When the cell is ready to divide and pass along the genetic code to a newly formed cell, the chromatin condenses into 46 gene-containing units, called chromosomes, which are grouped in 23 homologous pairs. One pair determines sex (the sex chromosomes), and the remaining 22 pairs (autosomes) determine the non-sex-related features of the progeny. In fertilization, 23 chromosomes (one-half of each pair) are donated by one parent, and 23 chromosomes are donated by the other parent; the zygote (fertilized egg) then contains 23 chromosome pairs (46 chromosomes).

The effect of a mutant gene inherited from one parent depends on whether the other parent donates a similar mutant gene (homozygous) or a normal gene (heterozygous). In general, heterozygotes (those with only one mutant gene) have about 50 percent of the normal enzyme activity because only one gene is normal, although this may be sufficient to maintain good health. When, however, the gene directs the synthesis of an enzyme that is rate limiting (*i.e.*, a series of reactions will proceed at a normal rate only as long as this enzyme is available), then even a 50 percent reduction in enzyme activity may produce serious disease. This gene and the disease that results are called autosomal dominant. Homozygotes (those with two mutant genes) suffer from more severe adverse effects

because the mutant enzyme has little or no normal activity. When two mutant genes must be paired to produce a disease, the gene and the disease are autosomal recessive. These basic points are summarized in Table 1.

If the mutant gene is part of the X chromosome (a sex chromosome; females have two X chromosomes per pair, and males have an X chromosome and a Y chromosome), the disease is sex-linked. All male offspring are affected because the Y chromosome of the XY pair does not have a compensating normal gene. Because the mutation is on the X chromosome, however, and males transmit only the Y chromosome to their sons during fertilization, males do not transmit the disease to them. They do, however, transmit the carrier state (*i.e.*, the mutant X chromosome) to their daughters. If the X-linked disease is dominant, both heterozygote and homozygote females are affected, the latter more severely; if the disease is recessive, only female homozygotes are affected. Female heterozygotes transmit an X chromosome with a mutant gene to half their male offspring, prompting a disease, and an X chromosome with a mutant gene to half of their female offspring, who may in turn transmit the disease to their sons. Occasionally, only one child in a family has a metabolic disorder and there is no evidence of an abnormal gene in either parent. The abnormal gene in the child must be assumed to have arisen by a spontaneous mutation.

Two different enzyme deficiencies, operating at separate steps in the same metabolic pathway and coded for by two entirely different genes, can produce similar diseases. Alternatively, more than one type of mutant gene (allele) may have arisen at the same gene location (locus) so that diseases may arise in heterozygotes who have two different mutant genes. Some of the generalized metabolic diseases, such as diabetes, may be polygenic (*i.e.*, the disease is not fully expressed unless two or more gene loci—even on different chromosomes—are affected). Finally, certain gene-related metabolic disorders may not be expressed without a particular type of diet, hormonal stimulus, or drug.

If an inborn error of metabolism is suspected, a physician can investigate at several levels: (1) a thorough family history is obtained to analyze the pattern of inheritance; (2) blood or urine is analyzed to determine whether the substrate and the product of the suspect biochemical reaction are at normal levels; (3) the activity of the suspect enzyme may be measured in red or white blood cells, in tissue samples, or in skin cells grown in test tube cultures; (4) the number of enzyme molecules may be measured in blood or tissue samples; (5) the DNA molecule that codes for the suspected abnormal enzyme may be isolated from blood cells, and its nucleotide base sequence determined and identified as mutant; and (6) testing also may be applied to parents, siblings, or offspring of the affected person. The diagnosis of a treatable condition may be made through screening at birth. A severe metabolic disorder may be diagnosed in the fetus by using samples from the amniotic fluid surrounding the embryo. The information is then available for counseling or sometimes for initiating treatment before birth.

There are a variety of treatments for metabolic disorders: (1) If the product of the reaction is deficient, it can be provided; (2) if the high level of substrate is toxic, the substrate can be diminished by restricting the diet, by activating its removal from the body, or by preventing its synthesis at an earlier step; (3) enzyme activity may be increased by administering a cofactor or certain drugs; and (4) the enzyme deficiency may be replaced with normal enzyme, although this approach has been successful only rarely. The ideal treatment would be to insert the normal gene itself into the cells of the patient so that normal amounts of the enzyme can be synthesized and the metabolic defect corrected directly. There are techniques for transplanting genes into animals, and they are being studied for humans.

DISORDERS OF CARBOHYDRATE METABOLISM

Glycogen storage diseases. The role of glucose as a critical source of energy for the brain and as the initial source of energy for muscles requires that it be stored in an easily releasable form. The storage form of glucose is

Chromosomes

Investigation of an inborn error

Table 1: Hereditary Aspects of Metabolic Disease

subject	gene 1	gene 2	enzyme activity	disease	parents
normal	normal	normal	100%	no	one, both, or neither may be heterozygote
heterozygote	normal	mutant	about 50%	no if recessive, yes if dominant	one must be at least heterozygote
homozygote	mutant	mutant	near 0%	yes	both must be at least heterozygote

a large, branched molecule known as glycogen. The liver serves as the glycogen (and therefore glucose) depot for the brain, while muscles maintain a glycogen supply of their own. Glycogen is synthesized from, and broken down to, glucose as shown in Figure 14.

Types

There are at least 10 well-recognized varieties of glycogen storage diseases, all of which are autosomal recessive disorders (Figure 14). In each, glycogen accumulates in the liver or muscles and sometimes in the kidneys. In type I glycogen storage disease, the last step in glucose release from the liver is defective, leading to low levels of glucose in the blood (hypoglycemia). While hypoglycemia may not bring about symptoms in some persons, it may ultimately produce coma or convulsions in others. If the rate of conversion of the immediate precursor of glucose, glucose 6-phosphate, to lactic acid increases, acidosis (greater acidity of the blood) results. A compensatory mechanism that diverts amino acids away from protein synthesis and toward glucose synthesis (gluconeogenesis) leads to stunted growth. Treatment of type I disease consists of frequent daytime feedings augmented by glucose and protein delivery at night. This regimen protects the brain from hypoglycemia, reduces the size of the liver, and restores growth.

Defects in earlier steps in glycogen breakdown in the liver cause glycogen storage disease types III, VI, VIII, IX, and X (Figure 14), which usually lead to milder versions of type I disease. Each responds to similar treatment. In type

IV disease, an abnormal unbranched glycogen molecule (amylopectin) accumulates and causes severe liver failure. In glycogen storage disease type II, an entirely separate pathway for degrading the glycogen present in cell storage organelles, called lysosomes, is defective, resulting in the accumulation of glycogen, in the heart, muscles, and liver; in severe cases this leads to early death from heart failure. In glycogen storage disease type V, the first step in muscle glycogen degradation is deficient, preventing the release of glucose. Strenuous exercise must be curtailed, and muscle cramps and even muscle damage can occur. In glycogen storage disease type VII, glycolysis (breakdown of glucose) is deficient in muscle, producing similar symptoms and an accumulation of glycogen.

Galactosemia. A classical autosomal recessive disorder, galactosemia usually is due to a defective component of the second major step in the metabolism of the sugar galactose (Figure 14). When galactose is ingested, as in milk, galactose 1-phosphate accumulates, causing severe liver damage as well as kidney abnormalities, mental retardation, and hypoglycemia. An alternative product of galactose metabolism, the sugar galactitol, is deposited in the lens of the eye, causing cataracts. A defect in the first step of galactose metabolism causes only the accumulation of galactose itself and cataracts. Early, complete, and permanent elimination of galactose from the diet prevents these adverse effects.

Fructose disorders. Fructosuria is an autosomal recessive disorder caused by a defect in the first step in the metabolism of the sugar fructose (Figure 14). Fructose is therefore excreted in the urine, but there are no other consequences and the condition is harmless. In contrast, in the autosomal recessive disorder hereditary fructose intolerance, a deficiency of a later key enzyme in fructose metabolism leads to the accumulation of fructose 1-phosphate, a toxic substance that causes vomiting and hypoglycemia initially and liver damage over a long period of time. Eliminating fructose, in the form of fruit or sucrose (table sugar), is effective treatment. Fructose 1,6-diphosphatase deficiency is an autosomal recessive disorder that prevents the synthesis of glucose from its precursors (Figure 14), causing hypoglycemia after a short period of fasting. Ingestion of fructose also causes hypoglycemia and lactic acidosis. Treatment consists of frequent feeding and avoiding fructose.

Hypo-glycemia

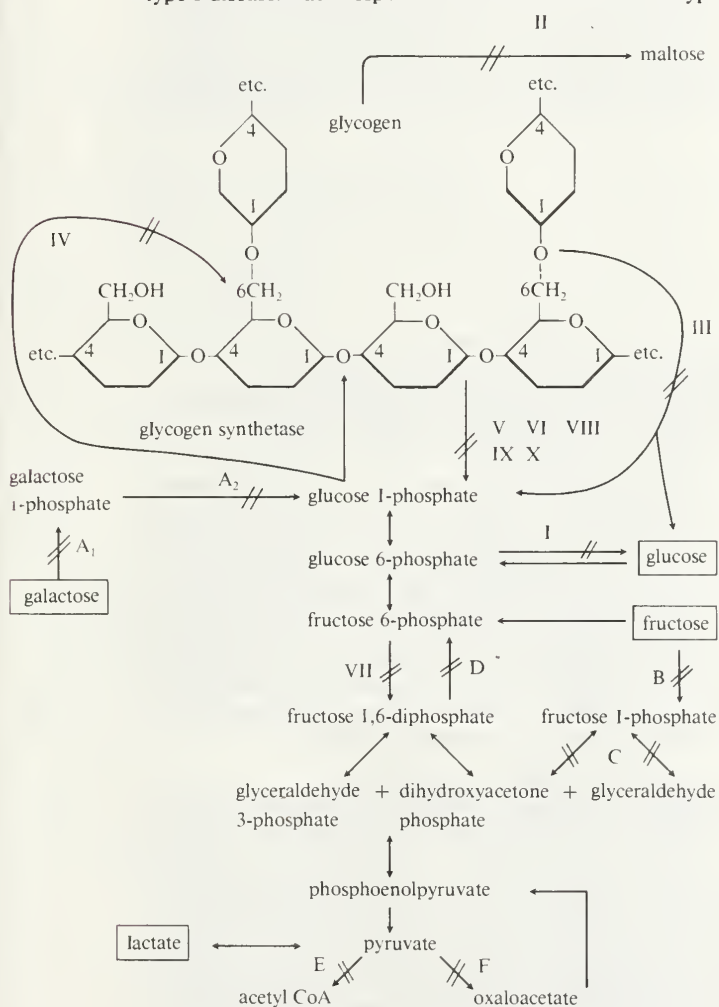


Figure 14: (I-X) Enzyme defects in glycogen storage diseases and (A₁-F) enzyme defects in other carbohydrate diseases. (I) Glucose-6-phosphatase (von Gierke). (II) Lysosomal α -1,4-glucosidase (Pompe). (III) Debranching enzyme (Forbes). (IV) Branching enzyme (Andersen). (V) Muscle phosphorylase (McArdle). (VI) Liver phosphorylase (Hers, VIII, IX, X). (VII) Muscle phosphofructokinase (Tauri). (A₁) Galactokinase. (A₂) Galactose 1-phosphate UDP transferase (galactosemia). (B) Fructokinase (fructosuria). (C) Aldolase (hereditary fructose intolerance). (D) Fructose 1,6-diphosphatase deficiency. (E) Pyruvate dehydrogenase complex deficiency. (F) Pyruvate carboxylase deficiency.

Pyruvate disorders. Two autosomal recessive disorders, pyruvate dehydrogenase complex deficiency and pyruvate carboxylase deficiency (Figure 14), can cause the accumulation of a three-carbon metabolite, pyruvate, in the blood. As a result, the level of lactic acid in the blood rises, causing acidosis. Both disorders also cause abnormalities of the central nervous system and retardation. In addition, pyruvate carboxylase deficiency causes hypoglycemia because glucose is not synthesized (Figure 14).

Intestinal carbohydrate malabsorption. Carbohydrates are ingested primarily in the form of larger molecules, such as starch (the chief storage form of carbohydrates in plants), from which the smaller sugars eventually are produced by the actions of intestinal enzymes (oligosaccharidases), such as maltase, which splits maltose to glucose; lactase, which splits lactose (milk sugar) to glucose and galactose; and sucrase, which splits sucrose (table sugar) to glucose and fructose. Lactase deficiency is common, occurring in up to 20 percent of North American whites, and 50 to 90 percent of blacks and Asians. Although hereditary, this disorder usually does not show symptoms until after infancy. Abdominal fullness, nausea, and cramping pain with watery diarrhea follow the intake of milk products containing lactose. Similar symptoms occur with a deficiency of sucrase after the ingestion of table sugar. In these conditions, glucose levels do not increase in the blood after lactose or sucrose, respectively, are ingested. These symptoms are eliminated when lactose or sucrose, respectively, are avoided.

Mucopolysaccharidoses. Sugar molecules may come together in repeating units to form polymers known as mucopolysaccharides. When proportionately more sugar polymers are bound to a protein core, the molecule is known as a proteoglycan. When the molecule contains pro-

portionately more protein than sugar polymers, it is called a glycoprotein. All of these perform diverse functions, particularly as basic components of the material between the tissue cells (the matrix). Mucopolysaccharides are regularly degraded as part of the process of tissue renewal. If any one of 15 enzymes responsible for this degradation is abnormal, partially degraded mucopolysaccharides accumulate in the tissues, causing recessive diseases with similar characteristics (Hurler, Hunter, Sanfilippo, and Morquio syndromes). These characteristics may include short stature, various deformities of the skeleton, stiffened joints, corneal clouding, deafness, dental abnormalities, hardening of the arteries, enlarged liver and spleen, and mental retardation. These syndromes are diagnosed by analyzing the incompletely degraded mucopolysaccharides in the urine and tissues. Prenatal diagnosis is also possible. No treatment is available for any of these disorders.

DISORDERS OF LIPID METABOLISM

Blood lipid disorders. Lipid molecules, such as triglycerides and cholesterol, circulate in the blood and are bound to carrier protein molecules. In certain hereditary metabolic disorders, the levels of triglycerides and cholesterol in the blood are high because the mechanisms that remove them are defective. As a result, triglycerides and cholesterol may be deposited in the skin, tendons, and walls of blood vessels, which can lead to coronary heart disease and early death.

Familial
dysbeta-
lipoproteinemia

In familial dysbetalipoproteinemia (type III hyperlipemia, autosomal recessive), a mutant carrier protein causes the accumulation of triglycerides and cholesterol. The disease is treated by restricting fat and cholesterol in the diet and with drugs that lower the level of lipids in the blood. In familial hypercholesterolemia (type II hyperlipemia, autosomal dominant), a protein receptor on the surface of all cells, which allows cholesterol to move into the cell, is either missing or defective. As a result, cholesterol is not moved into the cells, where mechanisms are in place to register the availability of cholesterol to metabolic processes. Under normal conditions, when enough cholesterol is registered, feedback mechanisms signal the cholesterol-synthesizing enzymes to cease cholesterol synthesis. In the disease state, the synthesizing enzymes are relieved of feedback inhibition, thus inducing the synthesis of still more cholesterol. Homozygotes (one per 1,000,000 persons) have much higher levels of cholesterol in the blood than heterozygotes (one per 500 persons) and may die before the age of 20. Restricting cholesterol in the diet and drug treatment sufficiently lower cholesterol in heterozygotes but not in homozygotes. Although relatively uncommon, familial hypercholesterolemia is an important model for the association between a high prevalence of coronary heart disease and modest elevations of blood cholesterol levels, which is often found in the West. In familial lipoprotein lipase deficiency (type I hyperlipemia, autosomal recessive), the enzyme responsible for transporting triglycerides into cells is deficient. Triglycerides build up in the blood, and the blood serum takes on a creamy consistency and colour after the ingestion of fat. There is abdominal pain, inflammation of the pancreas, and enlargement of the liver, but there is no increased risk of coronary heart disease. Restricting dietary fat eliminates the symptoms.

In familial betalipoprotein deficiency (autosomal recessive), the absence of a major transport protein leads to low levels of triglyceride and cholesterol in the blood. Dietary fat and lipid-soluble vitamins, particularly vitamin E, are not adequately absorbed. Vitamin E deficiency may be partly responsible for the resulting severe dysfunction of the nervous system, the eyes, and the muscles.

In familial high-density lipoprotein deficiency (Tangier's disease, autosomal recessive), the level of cholesterol in the blood is low, while that of triglycerides is high or normal. Cholesterol is deposited in many tissues, including those of the nerves, cornea, lymph nodes, and enlarged yellow tonsils. There is no treatment.

Wolman's
disease

In Wolman's disease, cholesterol esters accumulate in many tissues because the specific enzymes that cleave them are absent. Massive enlargement of the liver and

spleen, calcification of the adrenal glands, failure to grow, and death in infancy are characteristic of this disease.

Lipid oxidative disorders. The oxidation of fatty acids for energy requires that they be transferred into the mitochondria by associating with a cofactor, carnitine, which is synthesized in the body. Deficiency of carnitine or of the enzyme carnitine acyltransferase (which facilitates the movement of fatty acids across the mitochondrial membrane) in muscle tissue leads to weakness and cramping with exercise. A low level of carnitine in the blood and tissues produces severe brain, liver, and heart damage. The exact defect leading to this autosomal recessive disorder is not known. Treatment with carnitine is partially effective.

Tissue lipid disorders. Complex molecules composed of long-chain fatty acids linked with various sugars, phosphate, and sulfate groups form essential components of cell membranes and are vital to the normal functioning of the nervous system. These molecules are degraded and resynthesized in the normal course of tissue turnover. Defects in the enzymes that regulate this process cause these lipid molecules to accumulate, which is associated with a number of severe autosomal recessive disorders (Figure 15).

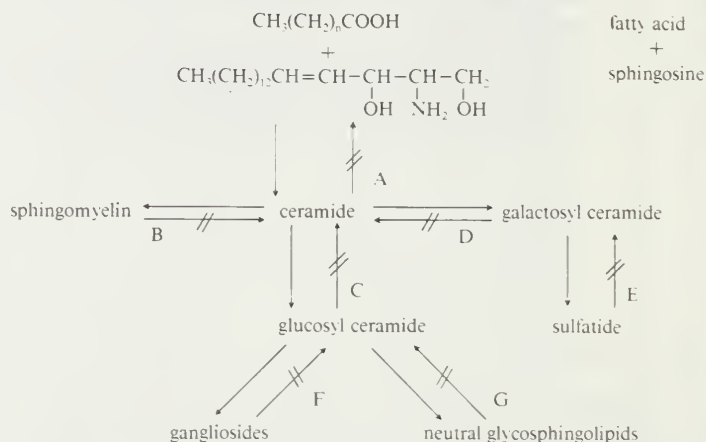


Figure 15: Enzyme defects in tissue lipid disorders. (A) Ceramidase (Farber). (B) Sphingomyelinase (Niemann-Pick). (C) Glucocerebrosidase (Gaucher). (D) Galactocerebroside β -galactosidase (Krabbe). (E) Arylsulfatase A (metachromatic leukodystrophy). (F) Hexosaminidase A (Tay-Sachs). (G) α -Galactosidase A (Fabry).

Among the symptoms are nervous system dysfunction and mental retardation; a unique cherry-red spot in the centre of the inner eye (Tay-Sachs disease); enlargement of the liver and spleen; characteristic foamy, lipid-filled cells in the bone marrow (Gaucher's and Niemann-Pick disease); or arthritis (Farber's disease). All but Gaucher's disease may cause early death. Fabry's disease is an X-linked disorder involving predominantly the blood vessels, kidneys, and eyes. For most of these conditions, the defective genes can be identified in carriers, and a prenatal diagnosis can be made. There is no successful treatment.

DISORDERS OF AMINO-ACID METABOLISM

Twenty amino acids, including nine that cannot be synthesized in humans and therefore must be obtained through food, are involved in metabolism. All are incorporated into proteins; some also are synthesized into important molecules in the body such as neurotransmitters, hormones, pigments, oxygen-carrying molecules, and the nucleotide bases of DNA and RNA. Each amino acid is further broken down to ammonia, carbon dioxide (CO_2), and water (H_2O). About 50 metabolic defects in amino-acid metabolism are known (Table 2), and their clinical consequences range from none to lethal.

Phenylalanine and tyrosine. The disorders involving the essential amino acids phenylalanine and tyrosine are examples of the many consequences of amino-acid dysfunction. Figure 16 demonstrates how these amino acids are synthesized and broken down into important molecules. Phenylketonuria is caused by the absence of the enzyme phenylalanine hydroxylase, which converts phenylalanine to tyrosine. Other variants (hyperphenylalaninemia IV and

Effects of phenylketonuria

V) are due to defects in the supply of the cofactor required for this reaction. Phenylketonuria produces severe mental retardation, seizures, rash, and high levels of phenylalanine in the blood and of its metabolite, phenylpyruvic acid, in the urine. Screening of blood at birth and restriction of phenylalanine in the diet are preventive measures.

Tyrosinemia results from a deficiency of the enzyme that catalyzes the first, and rate-limiting, step in tyrosine degradation, tyrosine transaminase (also called tyrosine aminotransferase). Tyrosinemia causes erosions of the skin and cornea of the eye. Deficiency in the next step causes a transient and harmless elevation of the level of tyrosine in the blood. A defect in the third successive enzymatic step causes an accumulation of homogentisic acid (which darkens the urine) and is called alkaptonuria. Although rare, alkaptonuria is of historical importance because it was Garrod's studies of this condition that led to his one-gene: one-enzyme hypothesis (one gene codes for one enzyme). When homogentisic acid is oxidized, the products are deposited in cartilage and connective tissue, producing a blue-gray colour in some areas of the skin and black discoloration in the joints; in adults, a severe form of arthritis, called ochronosis, results.

A disorder that results from a defect in the last step of tyrosine degradation, tyrosinosis, is a more generalized condition involving elevations of the levels of the amino acids methionine and tyrosine in the blood. Sudden liver failure may lead to early death; alternatively, there may be progressive kidney dysfunction in the form of Fanconi's syndrome (see below *Disorders of amino-acid transport*).

Tyrosine is a parent amino acid for skin, hair, and eye pigments (Figure 16). There are at least 10 syndromes, known generally as oculocutaneous albinism, that are

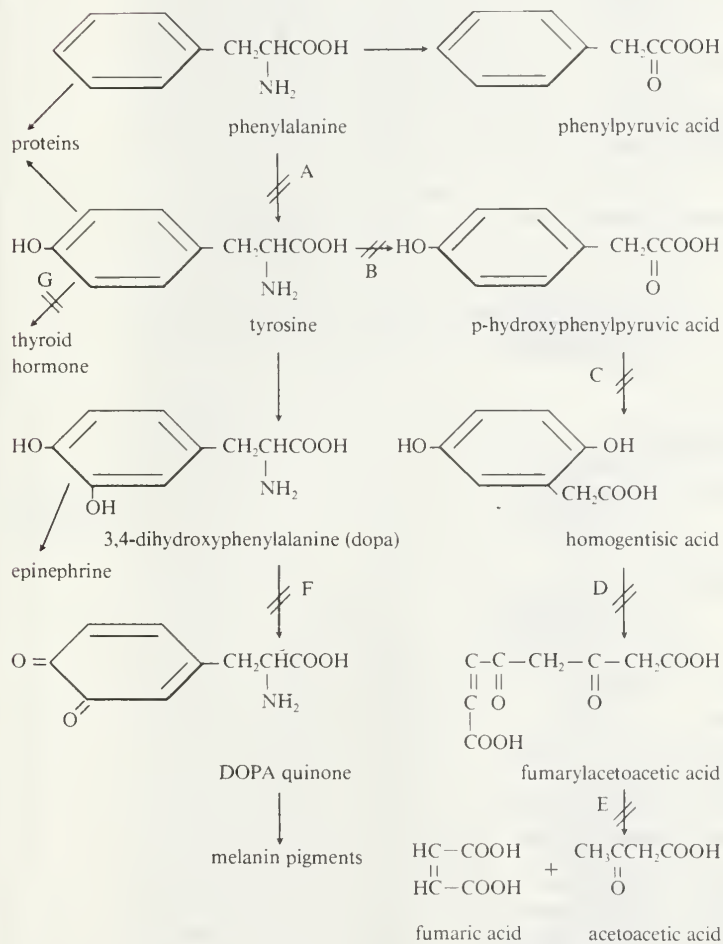


Figure 16: Enzyme defects in disorders of the amino acids phenylalanine and tyrosine.

(A) Phenylalanine hydroxylase (phenylketonuria). (B) Tyrosine transaminase (tyrosinemia). (C) p-Hydroxyphenylpyruvate oxidase (tyrosinemia). (D) Homogentisic acid oxidase (alkaptonuria). (E) Fumarylacetoacetate hydrolase (tyrosinosis). (F) Tyrosinase (albinism). (G) Peroxidase and others (familial goitre).

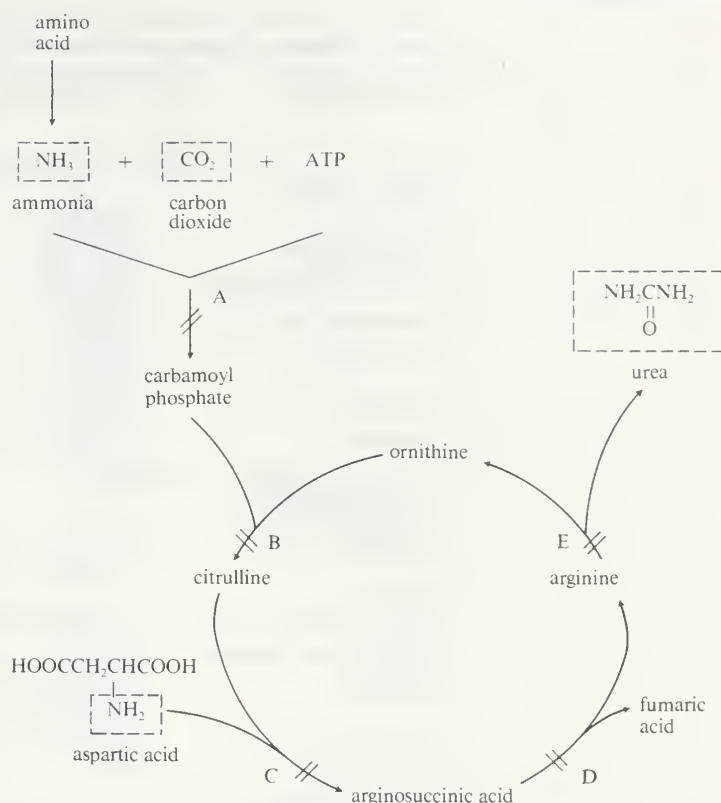


Figure 17: Enzyme defects in urea cycle disorders.

(A) Carbamoyl phosphate synthetase. (B) Ornithine carbamoyltransferase. (C) Argininosuccinate synthetase. (D) Argininosuccinate lyase. (E) Arginase.

characterized by the failure to form melanin pigments, resulting in partial or complete albinism. In addition to having white hair and white skin, persons with these defects are hypersensitive to light, and they have decreased vision and aberrant control of eye muscles. Oculocutaneous albinism can be autosomal recessive, autosomal dominant, or X-linked. A defect in the first step leading toward melanin synthesis, catalyzed by the enzyme tyrosinase, accounts for some of the syndromes; in others, the site of the block in pigment production is not known. The only available treatment is to avoid sunlight and to correct vision.

Finally, tyrosine is also the precursor amino acid for hormones of the thyroid gland. Defects in this synthetic pathway result in familial hypothyroidism, compensatory enlargement of the thyroid gland (goitre), severe growth failure, and retardation of central nervous system development. Screening of newborns and replacement of deficient thyroid hormone in affected persons prevent all or most of this damage.

Urea cycle enzymes. When amino acids are degraded, the ammonia nitrogen at one end of the molecule is split off, incorporated into urea (the chief solid component of urine), and excreted in the urine (Figure 17). This process starts when an ammonia group is combined with bicarbonate and phosphate to create carbamoyl phosphate, which hooks into a metabolic cycle involving five enzymes. The cycle takes in a second ammonia group from the amino acid aspartic acid, ultimately generating the amino acid arginine from which urea is split off. Defects at any of the first four steps of the cycle may cause a sudden increase in the level of ammonia in the blood, which in the newborn causes vomiting, coma, liver enlargement, and mental retardation. Early death can result when the enzyme is almost completely absent. A defect in the fifth step produces motor difficulties and seizures beginning in childhood. Specific amino-acid supplements, drugs that remove nitrogen by combining with the amino acids glycine and glutamine, and a diet low in protein are of limited benefit.

Branched chain amino acids. The amino acids valine, leucine, and isoleucine are degraded initially by removing

Formation of urea

Table 2: Disorders of Amino-Acid Metabolism

amino-acid disorder	enzyme deficiency	physical abnormalities	prenatal diagnosis feasible	early postnatal test available
Aromatic amino acids				
Phenylketonuria	phenylalanine hydroxylase	see text		X
Hyperphenylalanemia, types II and III	phenylalanine hydroxylase	see text		X
Hyperphenylalanemia, type IV	dihydropteridine reductase	see text		X
Hyperphenylalanemia, type V	dihydrobiopterin synthetase	see text		X
Tyrosinemia, type II	tyrosine transaminase	see text		X
Alkaptonuria	homogentisic acid oxidase	see text		X
Tyrosinosis	fumarylacetoacetate hydrolase	see text		
Oculocutaneous albinism, tyrosinase-negative type	tyrosinase	see text		X
Oculocutaneous albinism, tyrosinase-positive type S	unknown	see text		X
Familial goitre	peroxidase, deiodinase, and others	see text		X
Urea cycle				
Carbamoyl phosphate synthetase deficiency	carbamoyl phosphate synthetase	see text		X
Ornithine carbamoyltransferase* deficiency	ornithine carbamoyltransferase	see text	X	X
Citrullinemia	argininosuccinate synthetase	see text	X	X
Argininosuccinic acidemia	argininosuccinate lyase	see text	X	X
Hyperargininemia	arginase	see text		X
Branched chain amino acids				
Hypervalinemia	valine transaminase	see text	X	X
Hyperleucine-isoleucinemia	leucine-isoleucine transaminase	see text		X
Maple syrup urine disease (branched-chain ketoaciduria)	branched-chain α -keto acid decarboxylase	see text	X	X
Isovaleric acidemia	isovaleryl CoA dehydrogenase	see text	X	X
Glutaric aciduria, type II	acyl CoA dehydrogenases	see text		
3-Hydroxy-3-methylglutaric aciduria	3-hydroxy-3-methylglutaryl CoA lyase	see text	X	X
3-Ketothiolase deficiency	acetoacetyl CoA 3-ketothiolase	see text		
Propionic acidemia	propionyl CoA carboxylase	see text	X	X
Methylmalonic acidemia	methylmalonyl CoA mutase	see text	X	X
Others				
Histidinemia	histidase	speech disorder, mental retardation	X	X
5-Oxoprolinuria	glutathione synthetase	nervous system disorder, acidosis	X	X
Hyperprolinemia	proline oxidase or Δ^1 -pyrroline-5-carboxylate dehydrogenase	benign		X
Hyperhydroxyprolinemia	4-hydroxy-L-proline oxidase	benign		X
Prolidase deficiency (hyperimidodipeptiduria)	prolidase (imidodipeptidase)	rash, mental retardation		
Hyperornithinemia	ornithine δ -transaminase	visual disturbance, mental retardation	X	X
Hyperlysineemia, persistent	L-lysine: ketoglutarate reductase	uncertain	X	X
Homocystinuria	cystathionine β -synthetase	abnormal lens and bones, blood vessel clotting, mental retardation	X	X
Cystathioninuria	γ -cystathionase	benign	X	X
Cystinosis	unknown	kidney failure, cataracts	X	X
Hypersarcosinemia	sarcosine dehydrogenase complex	uncertain		X
Nonketotic hyperglycinemia	glycine cleavage enzyme complex	mental retardation, seizures		X
Primary oxaluria	α -ketoglutarate-glyoxylate carboligase	kidney stones and failure		X
Hyper- β -alaninemia	β -alanine transaminase	nervous system dysfunction		X
Carnosinase deficiency	serum carnosinase	uncertain		X
Homocarnosinosis	brain homocarnosinase	uncertain		X
Dihydrofolate reductase deficiency	dihydrofolate reductase	anemia, nervous system dysfunction		
Methylenetetrahydrofolate reductase deficiency	methylenetetrahydrofolate reductase	nervous system dysfunction		

*X-linked, dominant.

the amino group and producing CO₂, followed by reduction. Once degraded, however, their paths diverge. There are at least 11 genetic defects in the shared and nonshared portions of these pathways (Table 2). The consequences often include acidosis, hypoglycemia, and developmental retardation. In four of these conditions, the presence of incompletely degraded products in the urine is associated with distinctive odours, similar to those of maple syrup and sweaty feet. The appropriate amino acid or acids are restricted to a level that relieves the symptoms but still permits growth. Alkali, in the form of sodium bicarbonate, is given when the acidity of the blood increases suddenly.

Disorders of amino-acid transport. Energy is required to move many amino acids from the intestinal tract into the blood or to reclaim them from the urine by special cells in the kidney. This transport of amino acids is mediated by specific carrier proteins whose structures, like those of enzymes, are determined by genes. One gene may contain the code for the carrier protein that transports several amino acids; in addition, a single carrier protein may transport the same amino acids in the kidney and the intestinal tract. Mutant proteins with decreased transport activities may prevent the absorption of dietary amino acids or cause their loss in the urine. For example, in cystinuria there is decreased intestinal absorption and increased urinary excretion of four amino acids (cystine, arginine, lysine, and ornithine). The only consequence, however, is the formation of kidney stones. The essential amino acid lysine is probably absorbed from the diet in other forms. Iminoglycinuria is a completely harmless disorder of the absorption tubules in the kidney; glycine, proline, and hydroxyproline are excessively excreted in the urine without apparent consequence. In Hartnup disease, 13 amino acids, six of them essential, are not adequately reabsorbed from the intestine and kidney because of defects in the mechanism that removes them. As a result, high levels of these amino acids are excreted in the urine. Because the levels of the amino acids are low and because they are converted to toxic products within the intestinal tract, persons with Hartnup disease suffer intermittently from rashes and neuropsychiatric disturbances. In Fanconi's syndrome, because amino acids as well as other substrates (glucose, phosphate, and minerals, for example) are not absorbed well from the kidney, they are excreted in the urine. The syndrome occurs with several metabolic diseases, including cystinosis (a disorder in which the amino acid cystine is stored in cells of most body tissues). Bone disease, disturbed growth, dehydration, potassium deficiency, and acidosis result.

DISORDERS OF PORPHYRIN METABOLISM

Porphyryns are fluorescent cyclic molecules composed of four carbon-nitrogen rings that combine with specific metals, such as iron. These molecules, which carry oxygen in the blood (as hemoglobin) and in tissues, are also components of respiratory enzymes. Porphyryns are synthesized and degraded in the metabolic process shown in Figure 18. Defects in this process cause diseases, generally known as porphyrias, which primarily affect the skin (there is blistering on exposure to sunlight), the nervous system, and, in rare instances, the blood (rapid destruction of red blood cells). They are distinguished by the appearance of fluorescent products in the tissues and urine.

A dramatic form of porphyria is acute intermittent porphyria (AIP), which is caused by a deficiency of porphobilinogen deaminase, also called uroporphyrinogen I synthetase (Figure 18), an enzyme that catalyzes the second step in the pathway of porphyrin biosynthesis. Symptoms include abdominal pain, vomiting, muscle weakness, and mental disturbance. Porphobilinogen and its precursor, aminolevulinic acid, are found in the urine. AIP is an autosomal dominant disorder that occurs intermittently because the preceding reaction, catalyzed by the enzyme aminolevulinic acid synthetase, is the rate-limiting step in the pathway. Therefore, porphobilinogen does not reach toxic levels unless the activity of aminolevulinic acid synthetase is greatly increased by exposure to drugs, by starvation or excessive dieting, or by cyclic hormonal changes.

Red blood cells release hemoglobin (oxygen-carrying por-

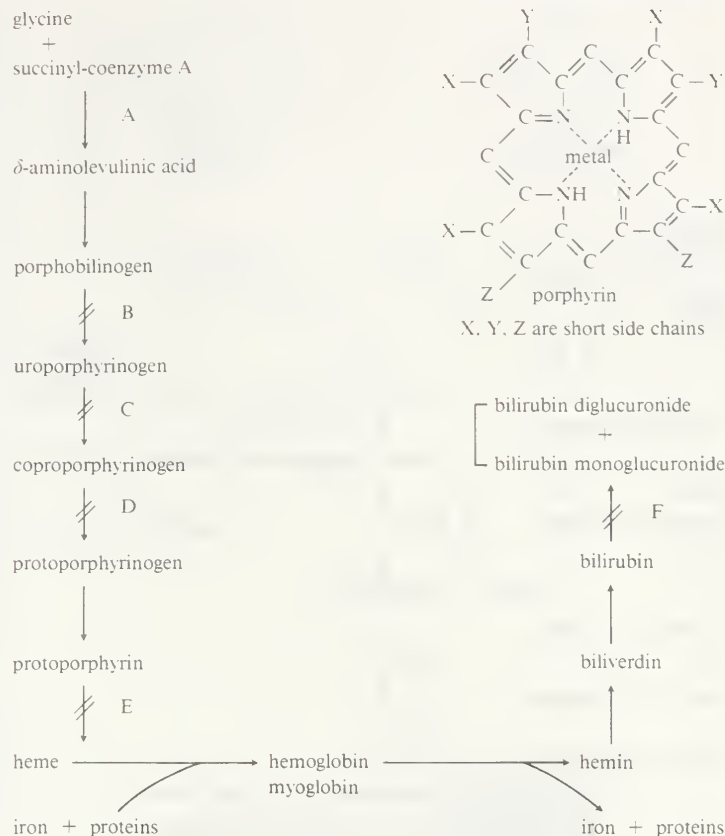


Figure 18: Enzyme defects in porphyrin metabolism. (A) δ -Aminolevulinic acid synthase. (B) Porphobilinogen deaminase, also called uroporphyrinogen I synthase (acute intermittent porphyria). (C) Uroporphyrinogen decarboxylase (porphyria cutanea tarda). (D) Coproporphyrinogen oxidase (hereditary coproporphyria). (E) Ferrochelatase (erythropoietic protoporphyria). (F) Bilirubin-UDP glucuronyltransferase (Crigler-Najjar and Gilbert).

phyrin molecules) when they die. The free porphyrins are degraded to a reddish yellow pigment known as bilirubin (Figure 18), which is transported to the liver. Defects in the mechanisms that remove and metabolize bilirubin increase the level of bilirubin in the blood and produce jaundice, a yellowing of the skin and eyes. A harmless defect in bilirubin metabolism is often seen for a short time in normal newborns and permanently in some adults. In infants without the enzyme bilirubin-UDP glucuronyl transferase (Crigler-Najjar syndrome), deposits of large amounts of bilirubin in the brain can cause severe and permanent damage to the nervous system. Treatment involves increasing enzyme activity by using the drug phenobarbital or converting bilirubin to less toxic products by exposing the infant to blue light.

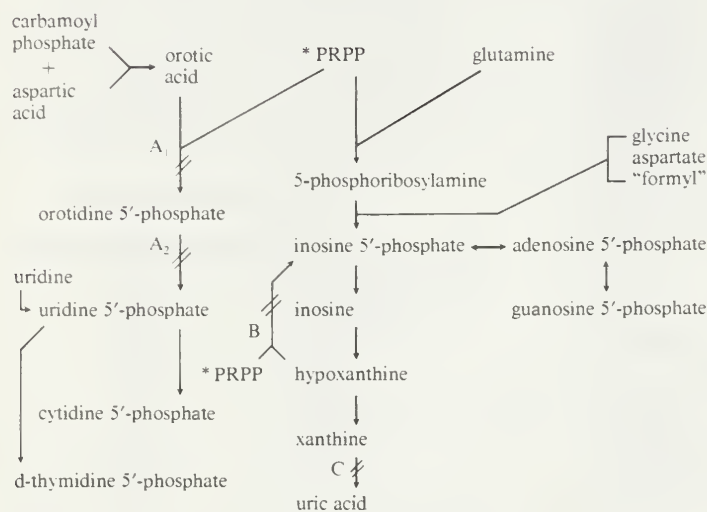
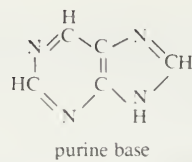
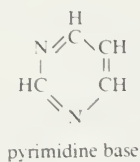
DISORDERS OF PURINE AND PYRIMIDINE METABOLISM

Purines and pyrimidines are the nucleotide bases in DNA and RNA. They have important roles in energy transformations (for example, ATP in the TCA cycle) and in other metabolic processes (for example, uridine triphosphate [UTP] in glycogen synthesis). Figure 19 shows a simplified overview of pyrimidine and purine synthesis to the parent compounds, uridine monophosphate and inosine monophosphate, respectively, as well as the important degradative route for the latter. Orotic aciduria (autosomal recessive) results from a block in the last two steps of uridine monophosphate synthesis. It is unique as a metabolic disorder in that one gene appears to control the activity of two successive enzymatic steps. Arrested growth, mental retardation, and anemia result from the deficiency of uridine monophosphate. Treatment with the pyrimidine base uridine is effective because it can be converted to uridine monophosphate by an alternate pathway.

Lesch-Nyhan syndrome, an X-linked disorder, results from a deficiency of the enzyme hypoxanthine-guanine phosphoribosyltransferase, which reuses hypoxanthine, a

Porphyryn
functions

Nucleotide
bases



* PRPP=5-phosphoribose 1-pyrophosphate

Figure 19: Enzyme defects in disorders of pyrimidine and purine metabolism.

(A₁) Orotate phosphoribosyltransferase (orotic aciduria).

(A₂) Orotidine 5'-phosphate decarboxylase (orotic aciduria).

(B) Hypoxanthine phosphoribosyltransferase (Lesch-Nyhan).

(C) Xanthine oxidase (xanthinuria).

product of inosine breakdown, for purine synthesis (Figure 19). If the enzyme is absent, the amount of unused key intermediate 5-phosphoribosyl 1-pyrophosphate (PRPP) is high, and the synthesis of inosine monophosphate continues (Figure 19), resulting in more hypoxanthine and, ultimately, uric acid. The nervous system is affected, resulting in mental retardation, abnormal movements, seizures, repeated self-mutilation, and death, usually by the age of 30. If the enzyme is only partially absent, gout results from the excess uric acid that settles in tissues. Gout also is a common disorder of purine metabolism in adults of the West, but in most cases the exact biochemical cause is unknown. The common denominator of all gouty conditions is high levels of uric acid in the blood. This leads

to the formation of uric acid crystals in joint fluids, causing attacks of arthritis, and in the urine, causing kidney stones. The drug allopurinol, which inhibits the conversion of the immediate precursor, xanthine, to uric acid by inhibiting the enzyme xanthine oxidase (Figure 19), is effective in preventing these consequences. Hereditary xanthinuria results from a deficiency of xanthine oxidase, but is usually harmless or, at most, associated with kidney stones formed from xanthine. (S.Ge.)

BIBLIOGRAPHY

General works: ALBERT L. LEHNINGER, *Principles of Biochemistry* (1982); LUBERT STRYER, *Biochemistry*, 2nd ed. (1981); EARLENE BROWN CUNNINGHAM, *Biochemistry: Mechanisms of Metabolism* (1978); JAY TEPPERMAN and HELEN M. TEPPERMAN, *Metabolic and Endocrine Physiology: An Introductory Text*, 5th ed. (1987); and S. DAGLEY and DONALD E. NICHOLSON, *An Introduction to Metabolic Pathways* (1970).

Cell metabolism: JAMES DARNELL, HARVEY LODISH, and DAVID BALTIMORE, *Molecular Cell Biology* (1986); T.A.V. SUBRAMANIAN (ed.), *Cell Metabolism, Growth and Environment*, 2 vol. (1986); W. BARTLEY, H.L. KORNBERG, and J.R. QUAYLE (eds.), *Essays in Cell Metabolism* (1970); JOHN M. DIETSCHY, ANTONIO M. GOTTO, JR., and JOSEPH A. ONTOK, *Disturbances in Lipid and Lipoprotein Metabolism* (1978); J. FRANK HENDERSON and A.R.P. PATERSON, *Nucleotide Metabolism: An Introduction* (1973); and DAVID A. BENDER, *Amino Acid Metabolism*, 2nd ed. (1985).

Regulation of metabolism: DANIEL E. ATKINSON, *Cellular Energy Metabolism and Its Regulation* (1977); E.A. NEWSHOLME and C. START, *Regulation in Metabolism* (1977); RONALD G. THURMAN, FREDERICK C. KAUFFMAN, and KURT JUNGERMANN (eds.), *Regulation of Hepatic Metabolism* (1986); and CHARLES ZAPSALIS and R. ANDERLE BECK, *Food Chemistry and Nutritional Biochemistry* (1985).

Genetics and metabolism: JAMES D. WATSON, *Molecular Biology of the Gene*, 4th ed. (1987); VINCENT M. RICCARDI, *The Genetic Approach to Human Disease* (1977); PHILIP F. BENSON and ANTHONY H. FENSOM, *Genetic Biochemical Disorders* (1985); AUBREY MILUNSKY (ed.), *Genetic Disorders and the Fetus: Diagnosis, Prevention, and Treatment*, 2nd ed. (1986); FORRESTER COCKBURN and RICHARD GITZELMAN, *Inborn Errors of Metabolism in Humans* (1982); and DAVID J. GALTON, *Molecular Genetics of Common Metabolic Disease* (1985).

Metabolic disorders: ARCHIBALD E. GARROD, *Inborn Errors of Metabolism* (1909, reprinted with a suppl. by H. HARRIS, 1963), a classic; JOHN B. STANBURY et al. (eds.), *The Metabolic Basis of Inherited Disease*, 5th ed. (1983); JOHN W. HARE (ed.), *Signs and Symptoms in Endocrine and Metabolic Disorders* (1986); PHILIP FELIG et al. (eds.), *Endocrinology and Metabolism* (1981); and F. DICKENS, P.J. RANDLE, and W.J. WHELAN (eds.), *Carbohydrate Metabolism and Its Disorders*, 3 vol. (1968-81). Special topics are treated in HENRIK GALBO, *Hormonal and Metabolic Adaptation to Exercise* (1983); and CHARLES S. LIEBER (ed.), *Metabolic Aspects of Alcoholism* (1977). (Ha.Ko./S.Ge.)

BOSTON PUBLIC LIBRARY



3 9999 04738 048 8

Brighton Branch Library
40 Academy Hill Road
Brighton, MA 02135-3316



WITHDRAWN

No longer the property of the
Boston Public Library.
Sale of this material benefits the Library.

WITHDRAWN

No longer the property of the
Boston Public Library.
Sale of this material benefits the Library.



